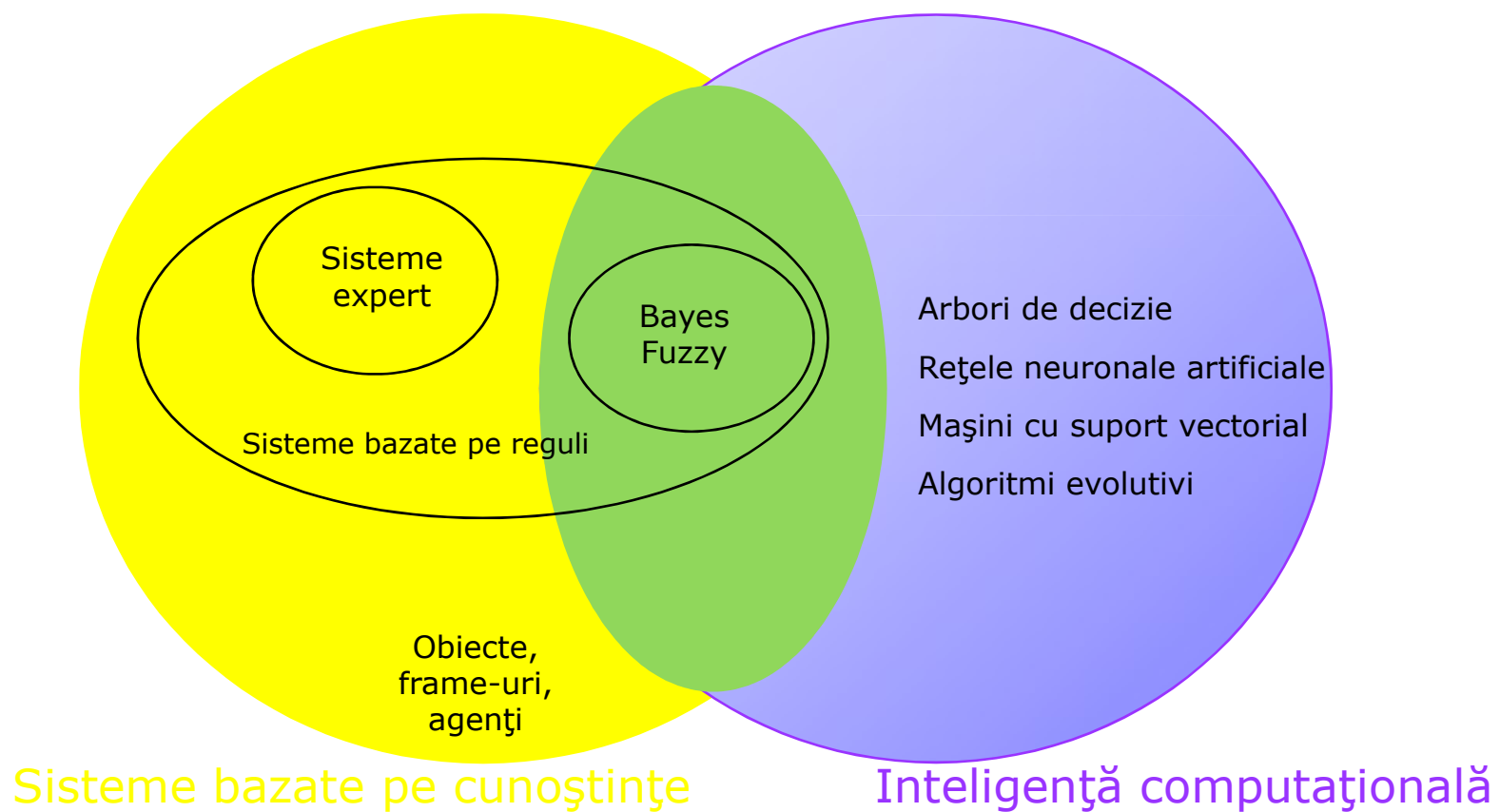


METODE INTELIGENTE DE REZOLVARE A PROBLEMELOR REALE



Laura Dioşan
Tema 4

Sisteme inteligente



Sisteme inteligente – SIS – Învățare automată

□ Tipologie

■ În funcție de experiența acumulată în timpul învățării:

- SI cu învățare supervizată
- SI cu învățare nesupervizată
- SI cu învățare activă
- SI cu învățare cu întărire

■ În funcție de modelul învățat (algoritmul de învățare):

- Arbori de decizie
- **Rețele neuronale artificiale**
- Algoritmi evolutivi
- Mașini cu suport vectorial
- Modele Markov ascunse

Sisteme inteligente – SIS – RNA

□ Rețele neuronale artificiale (RNA)

- Scop
- Definire
- Tipuri de probleme rezolvabile
- Caracteristici
- Exemplu
- Proiectare
- Evaluaire
- Tipologie

Sisteme inteligente – SIS – RNA

□ Scop

- Clasificare binară pentru orice fel de date de intrare (discrete sau continue)

- Datele pot fi separate de:

- o dreaptă $\rightarrow ax + by + c = 0$ (dacă $m = 2$)
- un plan $\rightarrow ax + by + cz + d = 0$ (dacă $m = 3$)
- un hiperplan $\sum a_i x_i + b = 0$ (dacă $m > 3$)

- Cum găsim valorile optime pt. a, b, c, d, a_i ?

- Rețele neuronale artificiale (RNA)
- Mașini cu suport vectorial (MSV)

- De ce RNA?

- Cum învață creierul?

Sisteme inteligente – SIS – RNA

□ Scop → De ce RNA?

- Unele sarcini pot fi efectuate foarte ușor de către oameni, însă sunt greu de codificat sub forma unor algoritmi
 - Recunoașterea formelor
 - vechi prieteni
 - caractere scrise de mână
 - vocea
 - Diferite raționamente
 - conducerea autovehiculelor
 - cântatul la pian
 - jucarea baschetului
 - înotul
- Astfel de sarcini sunt dificil de definit formal și este dificilă aplicarea unui proces de raționare pentru efectuarea lor

Sisteme inteligente – SIS – RNA

□ Scop → Cum învață creierul?

■ Creierul uman - componentă

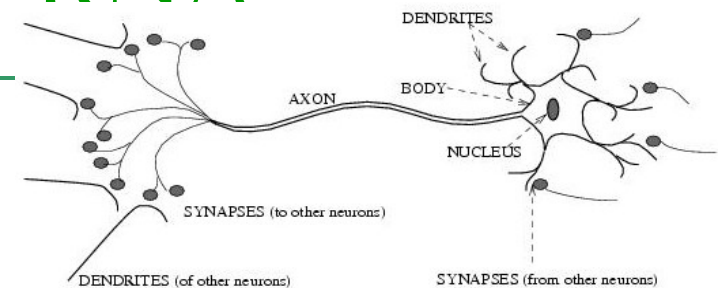
- Aproximativ 10.000.000.000 de neuroni conectați prin sinapse

□ Fiecare neuron

- are un corp (soma), un axon și multe dendrite
- poate fi într-una din 2 stări:
 - activ – dacă informația care intră în neuron depășește un anumit prag de stimulare –
 - pasiv – altfel

□ Sinapsă

- Legătura între axon-ul unui neuron și dendritele altui neuron
- Are rol în schimbul de informație dintre neuroni
 - 5.000 de conexiuni / neuron (în medie)
- În timpul vieții pot să apară noi conexiuni între neuroni



Sisteme inteligente – SIS – RNA

□ Scop → Cum învață creierul?

■ Cum "învață" (procesează informații)?

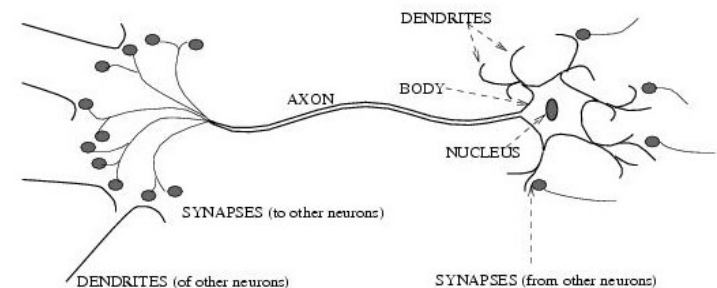
- Conexiunile care de-a lungul trăirii unor experiențe s-au dovedit utile devin permanente (restul sunt eliminate)
- Creierul este interesat de noutăți
- Modelul de procesare a informației
 - Învățare
 - Depozitare
 - Amintire

□ Memoria

- Tipologie
 - De scurtă durată
 - Imediată → 30 sec.
 - De lucru
 - De lungă durată
- Capacitate
 - Crește odată cu vârsta
 - Limitată → învățarea unei poezii pe strofe
- Influențată și de stările emoționale

□ Creierul

- rețea de neuroni
- sistem foarte complex, ne-liniar și paralel de procesare a informației
- Informația este depozitată și procesată de întreaga rețea, nu doar de o anumită parte a rețelei → informații și procesare globală
- Caracteristica fundamentală a unei rețele de neuroni → învățarea → rețele neuronale artificiale (RNA)



Sisteme inteligente – SIS – RNA

□ Definire

- Ce este o RNA?
- RN biologice vs. RN artificiale
- Cum învață rețeaua?

Sisteme inteligente – SIS – RNA

□ Definire → Ce este o RNA?

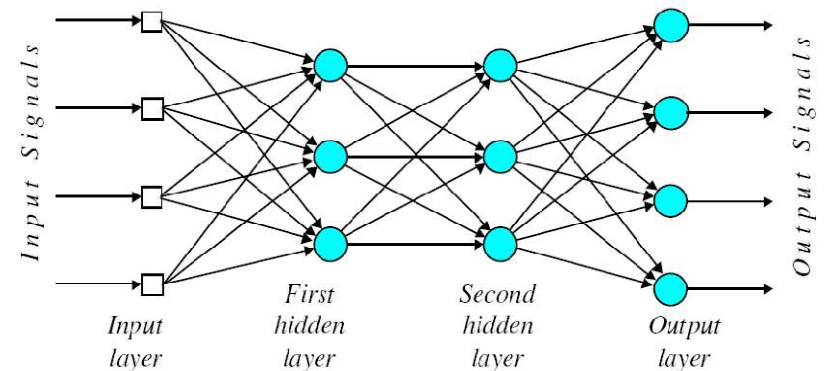
- O structură similară unei rețele neuronale biologice
- O mulțime de noduri (unități, neuroni, elemente de procesare) dispuse ca într-un graf pe mai multe straturi (*layere*)

□ Nodurile

- au intrări și ieșiri
- efectuează un calcul simplu prin intermediul unei funcții asociate → funcție de activare
- sunt conectate prin legături ponderate
 - Conexiunile între noduri conturează structura (arhitectura) rețelei
 - Conexiunile influențează calculele care se pot efectua

□ Straturile

- Strat de intrare
 - Conține m (nr de atribute al unei date) noduri
- Strat de ieșire
 - Conține r (nr de ieșiri) noduri
- Straturi intermediare (ascunse) – rol în “complicarea” rețelei
 - Diferite structuri
 - Diferite mărimi



Sisteme inteligente – SIS – RNA

- Definire → RN biologice vs. RN artificiale

| RNB | RNA |
|----------|---------------------|
| Soma | Nod |
| Dendrite | Intrare |
| Axon | Ieșire |
| Activare | Procesare |
| Synapsă | Conexiune ponderată |

- Definire → Cum învață rețeaua?

- Plecând de la un set de n date de antrenament de forma

$$((x_{p1}, x_{p2}, \dots, x_{pm}, y_{p1}, y_{p2}, \dots, y_{pr}))$$

cu $p = 1, 2, \dots, n$, m – nr atributelor, r – nr ieșirilor

- se formează o RNA cu m noduri de intrare, r noduri de ieșire și o anumită structură internă
 - un anumit nr de nivele ascunse, fiecare nivel cu un anumit nr de neuroni
 - cu legături ponderate între oricare 2 noduri
- se caută valorile optime ale ponderilor între oricare 2 noduri ale rețelei prin minimizarea erorii
 - diferența între rezultatul real y și cel calculat de către rețea

Sisteme inteligente – SIS – RNA

- Tipuri de probleme rezolvabile cu RNA
 - Datele problemei se pot reprezenta prin numeroase perechi atribut-valoare
 - Funcția obiectiv poate fi:
 - Unicriterială sau multicriterială
 - Discretă sau cu valori reale
 - Datele de antrenament pot conține erori (zgomot)
 - Timp de rezolvare (antrenare) prelungit

Sisteme inteligente – SIS – RNA

□ Proiectare

- Construirea RNA pentru rezolvarea problemei P
- Inițializarea parametrilor RNA
- Antrenarea RNA
- Testarea RNA

Sisteme inteligente – SIS – RNA

□ Proiectare

■ Construirea RNA pentru rezolvarea unei probleme P

- pp. o problemă de clasificare în care avem un set de date de forma:

- (x^d, t^d) , cu:
- $x^d \in \mathbf{R}^m \rightarrow x^d = (x^d_1, x^d_2, \dots, x^d_m)$
- $t^d \in \mathbf{R}^R \rightarrow t^d = (t^d_1, t^d_2, \dots, t^d_R)$,
- cu $d = 1, 2, \dots, n, n+1, n+2, \dots, N$

- primele n date vor fi folosite drept bază de antrenament a RNA
- ultimele $N-n$ date vor fi folosite drept bază de testare a RNA

- se construiește o RNA astfel:

- stratul de intrare conține exact m noduri (fiecare nod va citi una dintre proprietățile de intrare ale unei instanțe a problemei – $x^d_1, x^d_2, \dots, x^d_m$)
- stratul de ieșire poate conține R noduri (fiecare nod va furniza una dintre proprietățile de ieșire ale unei instanțe a problemei $t^d_1, t^d_2, \dots, t^d_R$)
- unul sau mai multe straturi ascunse cu unul sau mai mulți neuroni pe fiecare strat

Sisteme inteligente – SIS – RNA

□ Proiectare

- Construirea RNA pentru rezolvarea problemei P
- **Inițializarea parametrilor RNA**
- **Antrenarea RNA**
- Testarea RNA

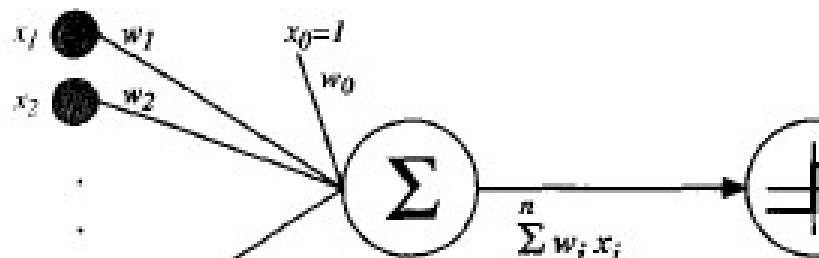
Sisteme inteligente – SIS – RNA

□ Proiectare

- Inițializarea parametrilor RNA
 - Inițializarea ponderile între oricare 2 noduri de pe straturi diferite
 - Stabilirea funcției de activare corespunzătoare fiecărui neuron (de pe straturile ascunse)
- Antrenarea (învățarea) RNA
 - Scop:
 - stabilirea valorii optime a ponderilor dintre 2 noduri
 - Algoritm
 - Se caută valorile optime ale ponderilor între oricare 2 noduri ale rețelei prin minimizarea erorii (diferența între rezultatul real y și cel calculat de către rețea)
 - Cum învață rețeaua?
 - Rețeaua = mulțime de unități primitive de calcul interconectate între ele →
 - Învățarea rețelei = \cup învățarea unităților primitive
 - Unități primitive de calcul
 - Perceptron
 - Unitate liniară
 - Unitate sigmoidală

Sisteme inteligente – SIS – RNA

- Proiectare → Antrenarea RNA → Cum învață rețeaua?
 - Neuronul ca element simplu de calcul
 - Structura neuronului
 - Fiecare nod are intrări și ieșiri
 - Fiecare nod efectuează un calcul simplu
 - Procesarea neuronului
 - Se transmite informația neuronului
 - Neuronul procesează informația
 - Se citește răspunsul neuronului
 - Învățarea neuronului – algoritmul de învățare a ponderilor care procesează corect informațiile
 - Se pornește cu un set inițial de ponderi oarecare
 - Cât timp nu este îndeplinită o condiție de oprire
 - Se procesează informația și se stabilește calitatea ponderilor curente
 - Se modifică ponderile astfel încât să se obțină rezultate mai bune



Sisteme inteligente – SIS – RNA

□ Proiectare → Antrenarea RNA → Cum învață rețeaua?

■ Neuronul ca element simplu de calcul

□ Structura neuronului

- Fiecare nod are intrări și ieșiri
- Fiecare nod efectuează un calcul simplu prin intermediul unei funcții asociate

□ Procesarea neuronului

- Se transmite informația neuronului → se calculează suma ponderată a intrărilor

$$net = \sum_{i=1}^n x_i w_i$$

- Neuronul procesează informația → se folosește o funcție de activare:

- Funcția constantă
- Funcția prag
- Funcția rampă
- Funcția liniară
- Funcția sigmoidală
- Funcția Gaussiană
- Funcția Relu

Sisteme inteligente – SIS – RNA

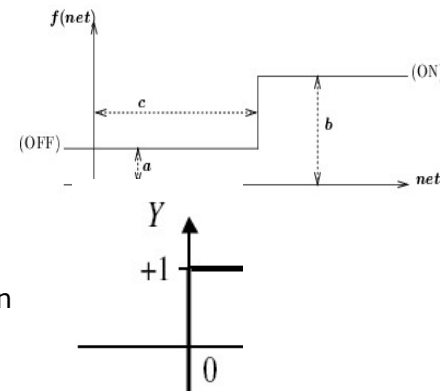
□ Proiectare → Antrenarea RNA → Cum învață rețeaua?

■ Funcția de activare a unui neuron

- Funcția constantă $f(net) = \text{const}$
- Funcția prag (c - pragul)

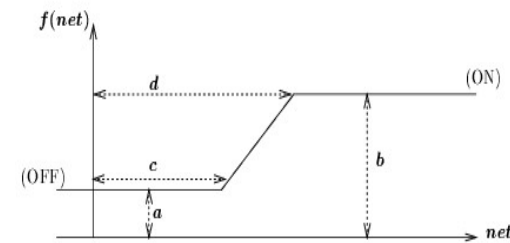
$$f(net) = \begin{cases} a, & \text{dacă } net < c \\ b, & \text{dacă } net > c \end{cases}$$

- Pentru $a=+1$, $b=-1$ și $c=0$ → funcția semn
- Funcție discontinuă



□ Funcția rampă

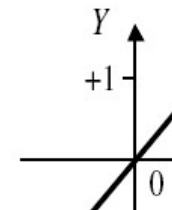
$$f(net) = \begin{cases} a, & \text{dacă } net \leq c \\ b, & \text{dacă } net \geq d \\ a + \frac{(net-c)(b-a)}{d-c}, & \text{altfel} \end{cases}$$



□ Funcția liniară

$$f(net) = a * net + b$$

- Pentru $a=1$ și $b=0$ → funcția identitate $f(net)=net$
- Funcție continuă



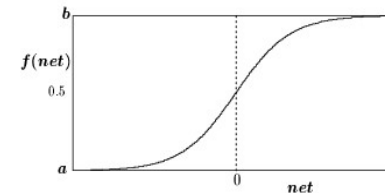
Sisteme inteligente – SIS – RNA

□ Proiectare → Antrenarea RNA → Cum învață rețeaua?

■ Funcția de activare a unui neuron

□ Funcția sigmoidală

- În formă de S
- Continuă și diferențiabilă în orice punct
- Simetrică rotațional față de un anumit punct ($net = c$)
- Atinge asimptotic puncte de saturație



$$\lim_{net \rightarrow -\infty} f(net) = a \quad \lim_{net \rightarrow \infty} f(net) = b$$

- Exemple de funcții sigmoidale:

$$f(net) = z + \frac{1}{1 + \exp(-x \cdot net + y)}$$

$$f(net) = \tanh(x \cdot net - y) + z$$

$$\text{unde } \tanh(u) = \frac{e^u - e^{-u}}{e^u + e^{-u}}$$

- Pentru $y=0$ și $z = 0 \rightarrow a=0, b = 1, c=0$
- Pentru $y=0$ și $z = -0.5 \rightarrow a=-0.5, b = 0.5, c=0$
- Cu cât x este mai mare, cu atât curba este mai abruptă

Sisteme inteligente – SIS – RNA

□ Proiectare → Antrenarea RNA → Cum învață rețeaua?

■ Funcția de activare a unui neuron

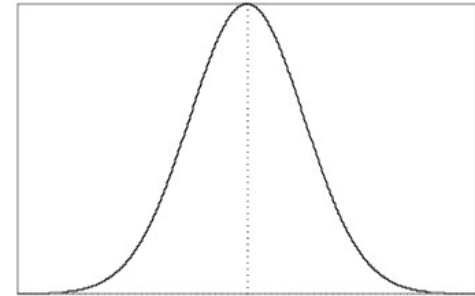
□ Funcția Gaussiană

- În formă de clopot
- Continuă
- Atinge asimptotic un punct de saturație

$$\lim_{net \rightarrow \infty} f(net) = a$$

- Are un singur punct de optim (maxim) – atins când $net = \mu$
- Exemplu

$$f(net) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\left(\frac{net - \mu}{\sigma}\right)^2\right]$$



Sisteme inteligente – SIS – RNA

□ Proiectare → Antrenarea RNA → Cum învață rețeaua?

■ Funcția de activare a unui neuron

□ Funcția ReLU

- În formă de rampă
- Continuă, monotonă
- Derivata ei este monotonă
- Codomeniu pozitiv $[0, \infty)$

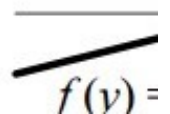
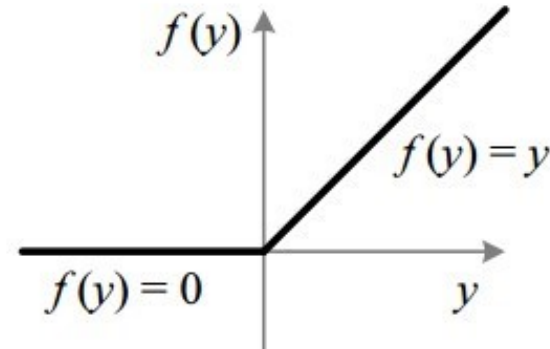
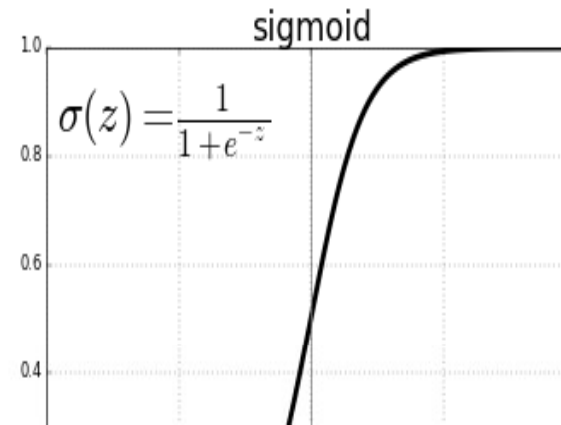
$$f(net) = \max(0, net)$$

$$f(net) = \begin{cases} 0, & \text{dacă } net < 0 \\ net, & \text{dacă } net \geq 0 \end{cases}$$

■ Variantă: Leaky ReLU

- Compensează problemele cu argumentele negative dint ReLU

$$f(net) = \begin{cases} a \cdot net, & \text{dacă } net < 0 \\ net, & \text{dacă } net \geq 0 \end{cases}$$



Sisteme inteligente – SIS – RNA

□ Proiectare → Antrenarea RNA → Cum învață rețeaua?

■ Neuronul ca element simplu de calcul

□ Structura neuronului

□ Procesarea neuronului

- Se transmite informația neuronului → se calculează suma ponderată a intrărilor

$$net = \sum_{i=1}^n x_i w_i$$

- Neuronul procesează informația → se folosește o funcție de activare:

- Funcția constantă
- Funcția prag
- Funcția rampă
- Funcția liniară
- Funcția sigmoidală
- Funcția Gaussiană

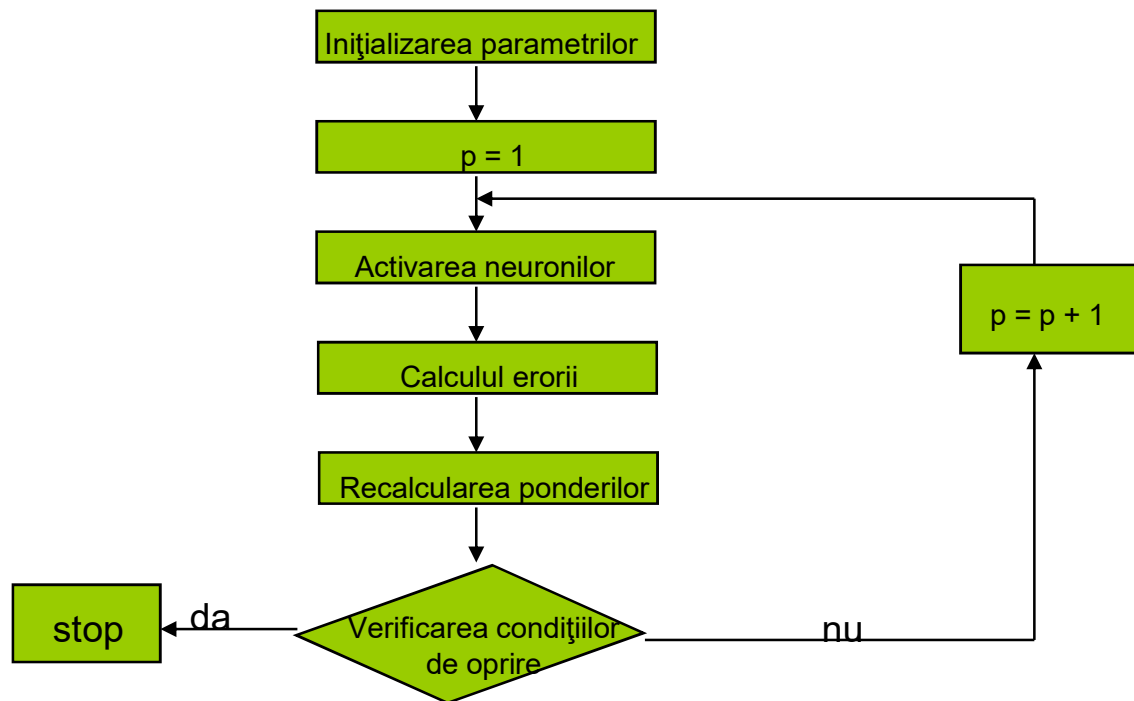
- Se citește răspunsul neuronului → se stabilește dacă rezultatul furnizat de neuron coincide sau nu cu cel dorit (real)

Sisteme inteligente – SIS – RNA

□ Proiectare → Antrenarea RNA → Cum învață rețeaua?

■ Neuronul ca element simplu de calcul

- Structura neuronului
- Procesarea neuronului
- Învățarea neuronului
 - Algoritm



Sisteme inteligente – SIS – RNA

□ Proiectare → Antrenarea RNA → Cum învață RNA?

■ Învățarea neuronului

□ 2 reguli de bază

■ Regula perceptronului → algoritmul perceptronului

1. Se porneste cu un set de ponderi oarecare
2. Se stabilește calitatea modelului creat pe baza acestor ponderi pentru **UNA** dintre datele de intrare
3. Se ajustează ponderile în funcție de calitatea modelului
4. Se reia algoritmul de la pasul 2 până când se ajunge la calitate maximă

■ Regula Delta → algoritmul scăderii după gradient

1. Se porneste cu un set de ponderi oarecare
 2. Se stabilește calitatea modelului creat pe baza acestor ponderi pentru **TOATE** dintre datele de intrare
 3. Se ajustează ponderile în funcție de calitatea modelului
 4. Se reia algoritmul de la pasul 2 până când se ajunge la calitate maximă
- Similar regulii perceptronului, dar calitatea unui model se stabilește în funcție de toate datele de intrare (tot setul de antrenament)

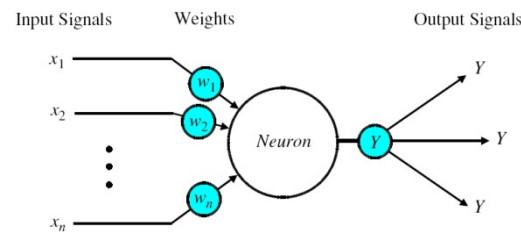
Sisteme inteligente – SIS – RNA

□ Proiectare → Antrenarea RNA → Cum învață RNA?

■ Învățarea neuronului

□ Pp că avem un set de date de antrenament de forma:

- (x^d, t^d) , cu:
 - $x^d \in \mathbf{R}^m \rightarrow x^d = (x^d_1, x^d_2, \dots, x^d_m)$
 - $t^d \in \mathbf{R}^R \rightarrow t^d = (t^d_1, t^d_2, \dots, t^d_R)$, și $R = 1$ (adică $t^d = (t^d_1)$)
 - cu $d = 1, 2, \dots, n$
- RNA = unitate primitivă de calcul (un neuron) → o rețea cu:
 - m noduri de intrare
 - legate de neuronul de calcul prin ponderile w_i , $i = 1, 2, \dots, m$ și
 - cu un nod de ieșire



Sisteme inteligente – SIS – RNA

□ Proiectare → Antrenarea RNA → Cum învață RNA?

■ Învățarea neuronului

□ Algoritmul perceptronului

- Se bazează pe minimizarea erorii asociată unei instanțe din setul de date de antrenament
- Modificarea ponderilor pe baza erorii asociate unei instanțe din setul de antrenament

Inițializare ponderi din rețea

$w_i = \text{random}(a,b)$, unde $i=1,2,\dots,m$

$d = 1$

Cât timp mai există exemple de antrenament clasificate incorect

Se activează neuronul și se calculează ieșirea

Perceptron → funcția de activare este funcția semn (funcție prag de tip discret, nediferențiable)

$$o^d = \text{sign}(\mathbf{w}\mathbf{x}) = \text{sign}\left(\sum_{i=1}^m w_i x_i\right)$$

Se stabilește ajustarea ponderilor $\Delta w_i = \eta(t^d - o^d)x_i^d$, unde $i = 1,2,\dots,m$

unde η - rată de învățare

Se ajustează ponderile $w'_i = w_i + \Delta w_i$

Dacă $d < n$ atunci $d++$

Altfel $d = 1$

SfCâtTimp

Sisteme inteligente – SIS – RNA

□ Proiectare → Antrenarea RNA → Cum învață RNA?

■ Învățarea neuronului

□ Algoritmul scădere după gradient

- Se bazează pe eroarea asociată întregului set de date de antrenament
- Modificarea ponderilor în direcția dată de cea mai abruptă pantă a reducerii erorii $E(\mathbf{w})$ pentru tot setul de antrenament

$$E(\mathbf{w}) = \frac{1}{2} \sum_{d=1}^n (t^d - o^d)^2$$

- Cum se determină cea mai abruptă pantă? → se derivează E în funcție de \mathbf{w} (se stabilește gradientul erorii E)

$$\nabla E(\mathbf{w}) = \left(\frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots, \frac{\partial E}{\partial w_m} \right)$$

- Gradientul erorii E se calculează în funcție de funcția de activare a neuronului (care trebuie să fie diferențiabilă, deci continuă)

- Funcția liniară
$$f(net) = \sum_{i=1}^m w_i x_i^d$$

- Funcția sigmoidală
$$f(net) = \frac{1}{1 + e^{-\mathbf{w}\mathbf{x}}} = \frac{1}{1 + e^{-\sum_{i=1}^m w_i x_i^d}}$$

- Cum se ajustează ponderile?

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}, \text{ unde } i = 1, 2, \dots, m$$

Sisteme inteligente – SIS – RNA

□ Proiectare → Antrenarea RNA → Cum învață RNA?

■ Învățarea neuronului

□ Algoritmul scădere după gradient → calcularea gradientului erorii

■ Funcția liniară

$$f(net) = \sum_{i=1}^m w_i x_i^d$$

$$\frac{\partial E}{\partial w_i} = \frac{\partial \frac{1}{2} \sum_{d=1}^n (t^d - o^d)^2}{\partial w_i} = \frac{1}{2} \sum_{d=1}^n \frac{\partial (t^d - o^d)^2}{\partial w_i} = \frac{1}{2} \sum_{d=1}^n 2(t^d - o^d) \frac{\partial (t^d - \mathbf{w}\mathbf{x}^d)}{\partial w_i}$$

$$\frac{\partial E}{\partial w_i} = \sum_{d=1}^n (t^d - o^d) \frac{\partial (t^d - w_1 x_1^d - w_2 x_2^d - \dots - w_m x_m^d)}{\partial w_i} = \sum_{d=1}^n (t^d - o^d) (-x_i^d)$$

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i} = \eta \sum_{d=1}^n (t^d - o^d) x_i^d$$

■ Funcție sigmoidală

$$f(net) = \frac{1}{1 + e^{-\mathbf{w}\mathbf{x}}} = \frac{1}{1 + e^{-\sum_{i=1}^m w_i x_i^d}} \quad y = s(z) = \frac{1}{1 + e^{-z}} \Rightarrow \frac{\partial s(z)}{\partial z} = s(z)(1 - s(z))$$

$$\frac{\partial E}{\partial w_i} = \frac{\partial \frac{1}{2} \sum_{d=1}^n (t^d - o^d)^2}{\partial w_i} = \frac{1}{2} \sum_{d=1}^n \frac{\partial (t^d - o^d)^2}{\partial w_i} = \frac{1}{2} \sum_{d=1}^n 2(t^d - o^d) \frac{\partial (t^d - sig(\mathbf{w}\mathbf{x}^d))}{\partial w_i} = \sum_{d=1}^n (t^d - o^d) (1 - o^d) o^d (-x_i^d)$$

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i} = \eta \sum_{d=1}^n (t^d - o^d) (1 - o^d) o^d x_i^d$$

Sisteme inteligente – SIS – RNA

- Proiectare → Antrenarea RNA → Cum învață RNA?
 - Învățarea neuronului
 - Algoritmul scădere după gradient (ASG)

| ASG simplu | ASG stocastic |
|---|--|
| <p>Inițializare ponderi din rețea $w_i = \text{random}(a,b)$, unde $i=1,2,\dots,m$</p> <p>$d = 1$</p> <p>Cât timp nu este îndeplinită condiția de oprire</p> <p>$\Delta w_i = 0$, unde $i=1,2,\dots,m$</p> <p>Pentru fiecare exemplu de antrenament (x^d, t^d), unde $d=1,2,\dots,n$</p> <p>Se activează neuronul și se calculează ieșirea o^d</p> <p>funcția de activare = funcția liniară $\rightarrow o^d = \mathbf{w}\mathbf{x}^d$</p> <p>funcția de activare = funcția sigmoid $\rightarrow o^d = \text{sig}(\mathbf{w}\mathbf{x}^d)$</p> <p>Pentru fiecare pondere w_i, unde $i = 1,2,\dots,m$</p> <p>Se stabilește ajustarea ponderii $\Delta w_i = \Delta w_i - \eta \frac{\partial E}{\partial w_i}$</p> <p>Pentru fiecare pondere w_i, unde $i = 1,2,\dots,m$</p> <p>Se ajustează fiecare pondere w_i $w_i = w_i + \Delta w_i$</p> | <p>Inițializare ponderi din rețea $w_i = \text{random}(a,b)$, unde $i=1,2,\dots,m$</p> <p>$d = 1$</p> <p>Cât timp nu este îndeplinită condiția de oprire</p> <p>$\Delta w_i = 0$, unde $i=1,2,\dots,m$</p> <p>Pentru fiecare exemplu de antrenament (x^d, t^d), unde $d=1,2,\dots,n$</p> <p>Se activează neuronul și se calculează ieșirea o^d</p> <p>funcția de activare = funcția liniară $\rightarrow o^d = \mathbf{w}\mathbf{x}^d$</p> <p>funcția de activare = funcția sigmoid $\rightarrow o^d = \text{sig}(\mathbf{w}\mathbf{x}^d)$</p> <p>Pentru fiecare pondere w_i, unde $i = 1,2,\dots,m$</p> <p>Se stabilește ajustarea ponderilor $\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$</p> <p>Se ajustează ponderea w_i $w_i = w_i + \Delta w_i$</p> |

Sisteme inteligente – SIS – RNA

- Proiectare → Antrenarea RNA → Cum învață RNA?
 - Învățarea neuronului

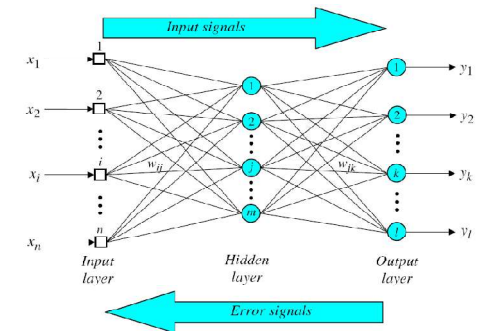
| Diferențe | Algoritmul perceptronului | Algoritmul scădere după gradient (regula Delta) |
|-----------------------------------|---|---|
| Ce reprezintă o^d | $o^d = \text{sign}(\mathbf{w}\mathbf{x}^d)$ | $o^d = \mathbf{w}\mathbf{x}^d$ sau $o^d = \text{sig}(\mathbf{w}\mathbf{x}^d)$ |
| Cum converge | Într-un nr finit de pași (până la separarea perfectă) | Asimptotic (spre eroarea minimă) |
| Ce fel de probleme se pot rezolva | Cu date liniar separabile | Cu orice fel de date (separabile liniar sau ne-liniar) |
| Ce tip de ieșire are neuronul | Discretă și cu prag | Continuă și fără prag |

Sisteme inteligente – SIS – RNA

□ Proiectare → Antrenarea RNA

■ Cum învață rețeaua?

- Rețeaua = mulțime de unități primitive de calcul interconectate între ele →
 - Învățarea rețelei = \cup învățarea unităților primitive
- Rețea cu mai mulți neuroni așezați pe unul sau mai multe straturi → RNA este capabilă să învețe un model mai complicat (nu doar liniar) de separare a datelor
- Algoritmul de învățare a ponderilor → backpropagation
 - Bazat pe algoritmul scădere după gradient
 - Îmbogățit cu:
 - Informația se propagă în RNA înainte (dinspre stratul de intrare spre cel de ieșire)
 - Eroarea se propagă în RNA înapoi (dinspre stratul de ieșire spre cel de intrare)



Se inițializează ponderile

Cât timp nu este îndeplinită condiția de oprire

Pentru fiecare exemplu (x^d, t^d)

Se activează fiecare neuron al rețelei

Se propagă informația înainte și se calculează ieșirea corespunzătoare fiecărui neuron al rețelei

Se ajustează ponderile

Se stabilește și se propagă eroarea înapoi

Se stabilesc erorile corespunzătoare neuronilor din stratul de ieșire

Se propagă aceste erori înapoi în toată rețeaua → se distribuie erorile pe toate conexiunile existente în rețea proporțional cu valorile ponderilor asociate acestor conexiuni

Se modifică ponderile

Sisteme inteligente – SIS – RNA

□ Proiectare → Antrenarea RNA

■ Cum învață o întreagă RNA?

□ Pp că avem un set de date de antrenament de forma:

■ (x^d, t^d) , cu:

- $x^d \in \mathbf{R}^m \rightarrow x^d = (x^d_1, x^d_2, \dots, x^d_m)$
- $t^d \in \mathbf{R}^R \rightarrow t^d = (t^d_1, t^d_2, \dots, t^d_R)$
- cu $d = 1, 2, \dots, n$

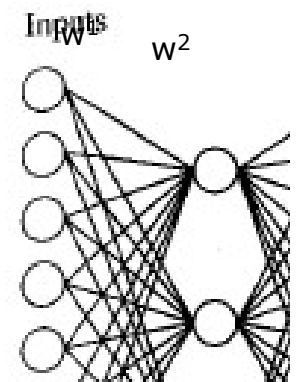
□ Presupunem 2 cazuri de RNA

■ O RNA cu un singur strat ascuns cu H neuroni → RNA₁

- m neuroni pe stratul de intrare,
- R neuroni pe stratul de ieșire,
- H neuroni pe stratul ascuns
- Ponderile între stratul de intrare și cel ascuns w_{ih}^1 cu $i=1,2,\dots,m$, $h=1,2,\dots,H$
- Ponderile între stratul ascuns și cel de ieșire w_{hr}^2 cu $h=1,2,\dots,H$ și $r=1,2,\dots,R$

■ O RNA cu p straturi ascunse, fiecare strat cu H_i ($i=1,2,\dots,p$) neuroni → RNA_p

- m neuroni pe stratul de intrare,
- R neuroni pe stratul de ieșire,
- P straturi ascunse
- H_p neuroni pe stratul ascuns p, $p=1,2,\dots,P$
- Ponderile între stratul de intrare și primul strat ascuns $w_{ih_1}^1$ cu $i=1,2,\dots,m$, $h_1=1,2,\dots,H_1$
- Ponderile între primul strat ascuns și cel de-al doilea strat ascuns $w_{h_1h_2}^2$ cu $h_1=1,2,\dots,H_1$ și $h_2=1,2,\dots,H_2$
- Ponderile între cel de-al doilea strat ascuns și cel de-al treilea strat ascuns $w_{h_2h_3}^3$ cu $h_2=1,2,\dots,H_2$ și $h_3=1,2,\dots,H_3$
- ...
- Ponderile între cel de-al p-1 strat ascuns și ultimul strat ascuns $w_{h_{p-1}h_p}^p$ cu $h_{p-1}=1,2,\dots,H_{p-1}$ și $h_p=1,2,\dots,H_p$
- Ponderile între ultimul strat ascuns și cel de ieșire $w_{h_p r}^{p+1}$ cu $h_p=1,2,\dots,H_p$ și $r=1,2,\dots,R$



Sisteme inteligente – SIS – RNA

- ❑ Proiectare → Antrenarea RNA → Cum învață o întreagă RNA?
 - Algoritmul backpropagation pentru RNA₁

Se inițializează ponderile w_{ih}^1 și w_{hr}^2 cu $i=1,2,\dots,m$, $h=1,2,\dots,H$ și $r=1,2,\dots,R$

Cât timp nu este îndeplinită condiția de oprire

Pentru fiecare exemplu (x^d, t^d)

Se activează fiecare neuron al rețelei

Se propagă informația înainte și se calculează ieșirea corespunzătoare fiecărui neuron al rețelei

$$o_h^d = \sum_{i=1}^m w_{ih}^1 x_i^d \text{ sau } o_h^d = \text{sig}\left(\sum_{i=1}^m w_{ih}^1 x_i^d\right), \text{ cu } h=1,2,\dots,H$$

$$o_r^d = \sum_{h=1}^H w_{hr}^2 o_h^d \text{ sau } o_r^d = \text{sig}\left(\sum_{h=1}^H w_{hr}^2 o_h^d\right), \text{ cu } r=1,2,\dots,R$$

Se ajustează ponderile

Se stabilește și se propagă eroarea înapoi

Se stabilesc erorile corespunzătoare neuronilor din stratul de ieșire

$$\delta_r^d = t_r^d - o_r^d \text{ sau } \delta_r^d = o_r^d (1 - o_r^d)(t_r^d - o_r^d), \text{ cu } r=1,2,\dots,R$$

Se modifică ponderile între nodurile de pe stratul ascuns și stratul de ieșire

$$w_{hr}^2 = w_{hr}^2 + \eta \delta_r^d o_h^d, \text{ unde } h=1,2,\dots,H \text{ și } r=1,2,\dots,R$$

Se propagă erorile nodurilor de pe stratul de ieșire înapoi în toată rețeaua → se distribuie erorile pe toate conexiunile existente în rețea proporțional cu valorile ponderilor asociate acestor conexiuni

$$\delta_h^d = \sum_{r=1}^R w_{hr}^2 \delta_r^d \text{ sau } \delta_h^d = o_h^d (1 - o_h^d) \sum_{r=1}^R w_{hr}^2 \delta_r^d$$

Se modifică ponderile între nodurile de pe stratul de intrare și stratul ascuns

$$w_{ih}^1 = w_{ih}^1 + \eta \delta_h^d x_i^d, \text{ unde } i=1,2,\dots,m \text{ și } h=1,2,\dots,H$$

Sisteme inteligente – SIS – RNA

- Proiectare → Antrenarea RNA → Cum învață o întreagă RNA?
 - Algoritmul backpropagation pentru RNA_p

Se inițializează ponderile $w_{ih_1}^1, w_{h_1h_2}^2, \dots, w_{h_{p-1}h_p}^p, w_{h_p r}^{p+1}$

Cât timp nu este îndeplinită condiția de oprire

Pentru fiecare exemplu (x^d, t^d)

Se activează fiecare neuron al rețelei

Se propagă informația înainte și se calculează ieșirea corespunzătoare fiecărui neuron al rețelei

$$o_{h_1}^d = \sum_{i=1}^m w_{ih_1}^1 x_i^d \text{ sau } o_{h_1}^d = \text{sig} \left(\sum_{i=1}^m w_{ih_1}^1 x_i^d \right), \text{ cu } h_1 = 1, 2, \dots, H_1$$

$$o_{h_2}^d = \sum_{h_1=1}^{H_1} w_{h_1h_2}^2 o_{h_1}^d \text{ sau } o_{h_2}^d = \text{sig} \left(\sum_{h_1=1}^{H_1} w_{h_1h_2}^2 o_{h_1}^d \right), \text{ cu } h_2 = 1, 2, \dots, H_2$$

...

$$o_{h_p}^d = \sum_{h_{p-1}=1}^{H_{p-1}} w_{h_{p-1}h_p}^p o_{h_{p-1}}^d \text{ sau } o_{h_p}^d = \text{sig} \left(\sum_{h_{p-1}=1}^{H_{p-1}} w_{h_{p-1}h_p}^p o_{h_{p-1}}^d \right), \text{ cu } h_p = 1, 2, \dots, H_p$$

$$o_r^d = \sum_{h_p=1}^{H_p} w_{h_p r}^{p+1} o_{h_p}^d \text{ sau } o_r^d = \text{sig} \left(\sum_{h_p=1}^{H_p} w_{h_p r}^{p+1} o_{h_p}^d \right), \text{ cu } r = 1, 2, \dots, R$$

Sisteme inteligente – SIS – RNA

- Proiectare → Antrenarea RNA → Cum învață o întreagă RNA?
 - Algoritmul backpropagation pentru RNA_p

Se inițializează ponderile $w_{ih_1}^1, w_{h_1h_2}^2, \dots, w_{h_{p-1}h_p}^p, w_{h_p r}^{p+1}$

Cât timp nu este îndeplinită condiția de oprire

Pentru fiecare exemplu (x^d, t^d)

Se activează fiecare neuron al rețelei

Se ajustează ponderile

Se stabilește și se propagă eroarea înapoi

Se stabilesc erorile corespunzătoare neuronilor din stratul de ieșire

$\delta_r^d = t_r^d - o_r^d$ sau $\delta_r^d = o_r^d(1 - o_r^d)(t_r^d - o_r^d)$, cu $r = 1, 2, \dots, R$
Se modifică ponderile între nodurile de pe ultimul strat ascuns și stratul de ieșire

$$w_{h_p r}^{p+1} = w_{h_p r}^{p+1} + \eta \delta_r^d o_{h_p}^d, \text{ unde } h_p = 1, 2, \dots, H_p \text{ și } r = 1, 2, \dots, R$$

Sisteme inteligente – SIS – RNA

- Proiectare → Antrenarea RNA → Cum învață o întreagă RNA?
 - Algoritmul backpropagation pentru RNA_p

Se inițializează ponderile $w_{ih_1}^1, w_{h_1h_2}^2, \dots, w_{h_{p-1}h_p}^p, w_{h_p r}^{p+1}$

Cât timp nu este îndeplinită condiția de oprire

Pentru fiecare exemplu (x^d, t^d)

Se activează fiecare neuron al rețelei

Se ajustează ponderile

Se stabilește și se propagă eroarea înapoi

Se stabilesc erorile corespunzătoare neuronilor din stratul de ieșire

Se modifică ponderile între nodurile de pe ultimul strat ascuns și stratul de ieșire

Se propagă (pe starturi) aceste erori înapoi în toată rețeaua → se distribuie erorile pe toate conexiunile existente în rețea proporțional cu valorile ponderilor asociate acestor conexiuni și se modifică ponderile corespunzătoare

$$\delta_{h_p}^d = \sum_{r=1}^R w_{h_p r}^{p+1} \delta_r^d \text{ sau } \delta_{h_p}^d = o_{h_p}^d (1 - o_{h_p}^d) \sum_{r=1}^R w_{h_p r}^{p+1} \delta_r^d$$

$$w_{h_p r}^{p+1} = w_{h_p r}^{p+1} + \eta \delta_r^d o_{h_p}^d, \text{ unde } h_p = 1, 2, \dots, H_p \text{ și } r = 1, 2, \dots, R$$

$$\delta_{h_{p-1}}^d = \sum_{h_p=1}^{H_p} w_{h_{p-1}h_p}^p \delta_{h_p}^d \text{ sau } \delta_{h_{p-1}}^d = o_{h_{p-1}}^d (1 - o_{h_{p-1}}^d) \sum_{h_p=1}^{H_p} w_{h_{p-1}h_p}^p \delta_{h_p}^d$$

$$w_{h_{p-1}h_p}^p = w_{h_{p-1}h_p}^p + \eta \delta_{h_p}^d o_{h_{p-1}}^d, \text{ unde } h_{p-1} = 1, 2, \dots, H_{p-1} \text{ și } h_p = 1, 2, \dots, H_p$$

...

$$\delta_{h_1}^d = \sum_{h_2=1}^{H_2} w_{h_1h_2}^2 \delta_{h_2}^d \text{ sau } \delta_{h_1}^d = o_{h_1}^d (1 - o_{h_1}^d) \sum_{h_2=1}^{H_2} w_{h_1h_2}^2 \delta_{h_2}^d$$

$$w_{ih_1}^1 = w_{ih_1}^1 + \eta \delta_{h_1}^d x_i^d, \text{ unde } i = 1, 2, \dots, m \text{ și } h_1 = 1, 2, \dots, H_1$$

MIRPR - ANN

Sisteme inteligente – SIS – RNA

- Proiectare → Antrenarea RNA → Cum învață o întreagă RNA?
 - Algoritmul backpropagation
 - Condiții de oprire
 - S-a ajuns la eroare 0
 - S-au efectuat un anumit număr de iterații
 - La o iterație se procesează un singur exemplu
 - n iterații = o epocă

Sisteme inteligente – SIS – RNA

□ Proiectare

- Construirea RNA pentru rezolvarea problemei P
- Inițializarea parametrilor RNA
- Antrenarea RNA
- **Testarea RNA**

Sisteme inteligente – SIS – RNA

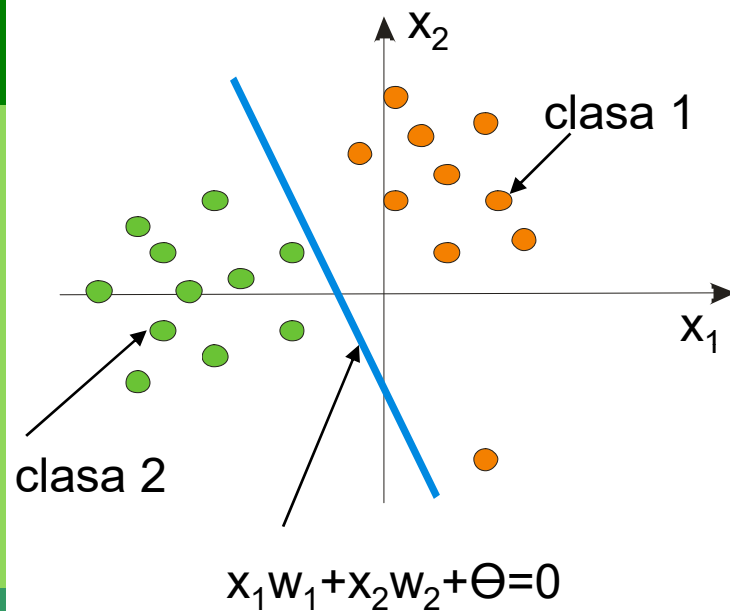
□ Proiectare

■ Testarea RNA

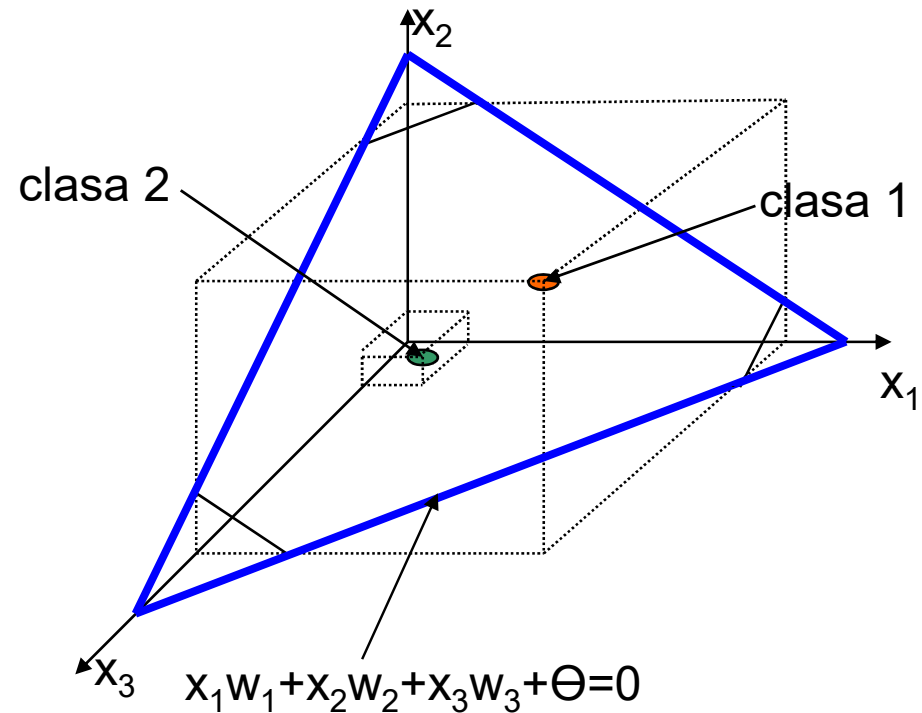
- Se decodifică modelul învățat de RNA
 - prin combinarea ponderilor cu intrările
 - ținând cont de funcțiile de activare a neuronilor și de structura rețelei

Sisteme inteligente – SIS – RNA

□ Exemplu



Clasificare binară cu $m=2$ intrări



Clasificare binară cu $m=3$ intrări

Sisteme inteligente – SIS – RNA

□ Exemplu

■ Perceptron pentru rezolvarea problemei *ȘI logic*

| Epoch | Inputs | | Desired output Y_d | Initial weights | | Actual output \hat{Y} | Error e | Final weights | |
|-------|--------|-------|-------------------------|-----------------|-------|----------------------------|--------------|---------------|-------|
| | x_1 | x_2 | | w_1 | w_2 | | | w_1 | w_2 |
| 1 | 0 | 0 | 0 | 0.3 | 0.1 | 0 | 0 | 0.3 | 0.1 |
| | 0 | 1 | 0 | 0.3 | -0.1 | 0 | 0 | 0.3 | -0.1 |
| | 1 | 0 | 0 | 0.3 | -0.1 | 1 | -1 | 0.2 | -0.1 |
| | 1 | 1 | 1 | 0.2 | 0.1 | 0 | 1 | 0.3 | 0.0 |
| 2 | 0 | 0 | 0 | 0.3 | 0.0 | 0 | 0 | 0.3 | 0.0 |
| | 0 | 1 | 0 | 0.3 | 0.0 | 0 | 0 | 0.3 | 0.0 |
| | 1 | 0 | 0 | 0.3 | 0.0 | 1 | -1 | 0.2 | 0.0 |
| | 1 | 1 | 1 | 0.2 | 0.0 | 1 | 0 | 0.2 | 0.0 |
| 3 | 0 | 0 | 0 | 0.2 | 0.0 | 0 | 0 | 0.2 | 0.0 |
| | 0 | 1 | 0 | 0.2 | 0.0 | 0 | 0 | 0.2 | 0.0 |
| | 1 | 0 | 0 | 0.2 | 0.0 | 1 | -1 | 0.1 | 0.0 |
| | 1 | 1 | 1 | 0.1 | 0.0 | 0 | 1 | 0.2 | 0.1 |
| 4 | 0 | 0 | 0 | 0.2 | 0.1 | 0 | 0 | 0.2 | 0.1 |
| | 0 | 1 | 0 | 0.2 | 0.1 | 0 | 0 | 0.2 | 0.1 |
| | 1 | 0 | 0 | 0.2 | 0.1 | 1 | -1 | 0.1 | 0.1 |
| | 1 | 1 | 1 | 0.1 | 0.1 | 1 | 0 | 0.1 | 0.1 |
| 5 | 0 | 0 | 0 | 0.1 | 0.1 | 0 | 0 | 0.1 | 0.1 |
| | 0 | 1 | 0 | 0.1 | 0.1 | 0 | 0 | 0.1 | 0.1 |
| | 1 | 0 | 0 | 0.1 | 0.1 | 0 | 0 | 0.1 | 0.1 |
| | 1 | 1 | 1 | 0.1 | 0.1 | 1 | 0 | 0.1 | 0.1 |

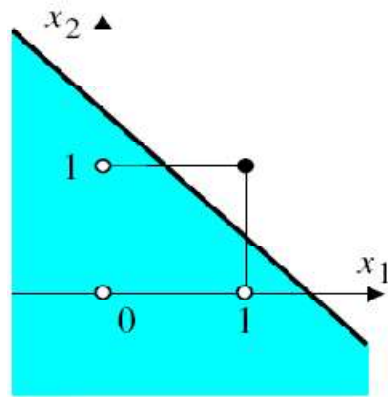
Threshold: $\theta = 0.2$; learning rate: $\alpha = 0.1$

Sisteme inteligente – SIS – RNA

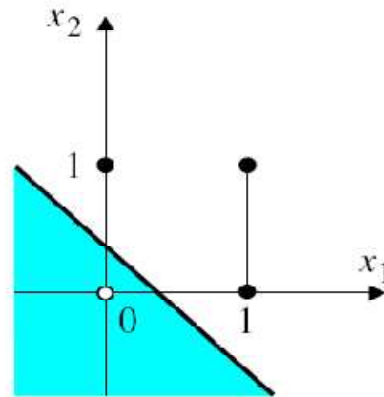
□ Exemplu

■ Perceptron - limitări

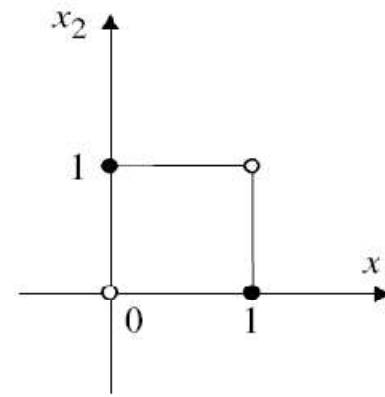
- Un perceptron poate învăța operațiile AND și OR, dar nu poate învăța operația XOR (nu e liniar separabilă)



(a) *AND* ($x_1 \cap x_2$)



(b) *OR* ($x_1 \cup x_2$)



(c) *Exclusive-OR*
($x_1 \oplus x_2$)

- Nu poate clasifica date non-liniar separabile

- soluții
 - Neuron cu un prag continuu
 - Mai mulți neuroni

Sisteme inteligente – SIS – RNA

□ Tipologie

■ RNA feed-forward

- Informația se procesează și circulă de pe un strat pe altul
- Conexiunile între noduri nu formează cicluri
- Se folosesc în special pentru învățarea supervizată
- Funcțiile de activare a nodurilor → liniare, sigmoidale, gaussiene

■ RNA recurente (cu feedback)

- Pot conține conexiuni între noduri de pe același strat
 - Conexiunile între noduri pot forma cicluri
 - RNA de tip Jordan
 - RNA de tip Elman
 - RNA de tip Hopfield
 - RNA auto-organizate → pentru învățarea nesupervizată
 - De tip Hebbian
 - De tip Kohonen (*Self organised maps*)
- } pentru învățarea supervizată

Sisteme inteligente – SIS – RNA

□ Avantaje

- Pot rezolva atât probleme de învățare super-vizată, cât și nesupervizată
- Pot identifica relații dinamice și neliniare între date
- Pot rezolva probleme de clasificare cu oricâte clase (multi-clasă)
- Se pot efectua calcule foarte rapid (în paralel și distribuit)

□ Dificultăți și limite

- RNA se confruntă cu problema overfitting-ului chiar și când modelul se învață prin validare încrucișată
- RNA pot găsi (uneori) doar optimele locale (fără să identifice optimul global)

Deep learning

□ Deep learning

- methodology in which we can train machine complex representations
- addresses the problem of learning hierarchical representations with a single (a few) algorithm(s)
- models with a feature hierarchy (lower-level features are learned at one layer of a model, and then those features are combined at the next level).
- it's deep if it has more than one stage of non-linear feature transformation
- Hierarchy of representations with increasing level of abstraction
 - Image recognition
 - Pixel → edge → texture → motif → part → object
 - Text
 - Character → word → word group → clause → sentence → story
 - Speech
 - Sample → spectral band → sound → ... → phone → phoneme → word

□ Deep networks/architectures

- Convolutional NNs
- Auto-encoders
- Deep Belief Nets (Restricted Boltzmann machines)
- Recurrent Neural Networks

Deep learning

- ❑ Collection of methods to improve the optimisation and generalisation of learning methods, especially NNs:
 - Rectified linear units
 - Dropout
 - Batch normalisation
 - Weight decay regularisation
 - Momentum learning
- ❑ Stacking layers of transformations to create successively more abstract levels of representation
 - Depth over breadth
 - Deep MLPs
- ❑ Shared parameters
 - Convolutional NNs
 - Recurrent NNs
- ❑ Technological improvements
 - Massively parallel processing: GPUs, CUDA
 - Fast libraries: Torch, cuDNN, CUDA-convNet, Theano

ML & optimisation

- ❑ An ML algorithm as an optimisation approach
 - An optimization problem
 - ❑ minimize the loss function
 - ❑ with respect to the parameters of the score function.
 - **score function**
 - ❑ maps the raw data to class scores/labels
 - **loss function**
 - ❑ quantifies the agreement between the predicted scores and the ground truth scores/labels
 - ❑ ANN: quantifies the quality of any particular set of weights **W**
 - ❑ two components
 - The data loss computes the compatibility between the computed scores and the true labels.
 - The regularization loss is only a function of the weights

Classification

□ Suppose a supervised classification problem

- Some input data (examples, instances, cases)
 - Training data – as pairs $(\text{attribute_data}_i, \text{label}_i)$, where
 - $i = 1, N$ ($N = \#$ of training data)
 - $\text{attribute_data}_i = (\text{atr}_{i1}, \text{atr}_{i2}, \dots, \text{atr}_{im})$, $m = \#$ attributes (characteristics, features) for an input data
 - $\text{label}_i \in \{\text{Label}_1, \text{Label}_2, \dots, \text{Label}_{\# \text{classes}}\}$
 - Test data – as $(\text{attribute_data}_i)$, $i = 1, n$ ($n = \#$ of testing data).
- Determine
 - An unknown function that maps inputs (features) into outputs (labels)
 - Output (label/class/value/score) associated to a new data by using the learnt function

□ Quality of learning

- Accuracy/Precision/Recall/etc
 - does not reflect the learnt decision model
- A loss function
 - Expresses (encodes) the learnt model
 - Difference between desired (D) and computed (C) output
 - L_2 norm - Quadratic cost (*mean squared error*) $\sum \|D - C\|^2$
 - L_1 norm $\sum |D - C|$
 - SVM loss (hinge loss, max-margin loss) $\sum_i \sum_{j, j \neq y_i} \max(C_j - D_{y_i} + \Delta, 0)$
 - Softmax loss $\sum_i [-\ln(\exp(D_{y_i}) / \sum_{j, j \neq y_i} \exp(C_j))]$
 - Cross-entropy $-\sum [D \ln C + (1 - D) \ln(1 - C)] / n$

Classifiers

□ Several important mappings

- Constant $f(x) = c$
- Step $f(x) = a$, if $x < \text{theta}$
 b , otherwise
- Linear $f(x) = a x + b$
- Sigmoid $\sigma(x) = 1/(1+e^{-x})$ (avoid it in a Conv NN)
- Hyperbolic tangent function $\tanh(x) = 2\sigma(2x) - 1$
- Rectified linear neuron/unit (ReLU) $f(x) = \max(0, x)$
- Leak ReLU (Parametric rectifier) $f(x) = \max(\alpha x, x)$
- Maxout $\max(w_1^T x + b_1, w_2^T x + b_2)$
- Exponential linear units (ELU) $f(x) = x$, if $x > 0$
 $\alpha (\exp(x) - 1)$, if $x \leq 0$

A linear classifier

$$\begin{aligned} f(x, w) &= w \cdot x + b, \\ w &\in \mathbb{R}^{\text{\#classes} \times \text{\#features}} \\ x &\in \mathbb{R}^{\text{\#features} \times 1} \\ b &\in \mathbb{R}^{\text{\#classes}} \end{aligned}$$

A non linear classifier

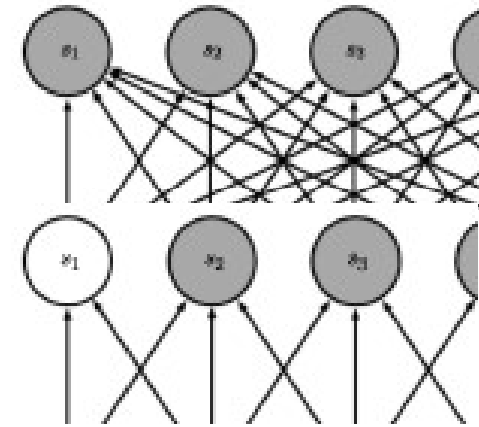
$$\begin{aligned} f(x, w) &= w_2 \max(0, w_1 \cdot x + b_1) + b_2, \\ w_1 &\in \mathbb{R}^{\text{PARAM} \times \text{\#features}} \\ x &\in \mathbb{R}^{\text{\#features} \times 1} \\ b_1 &\in \mathbb{R}^{\text{PARAM}} \\ w_2 &\in \mathbb{R}^{\text{\#classes} \times \text{PARAM}} \\ b_2 &\in \mathbb{R}^{\text{\#classes}} \end{aligned}$$

Classical ANN

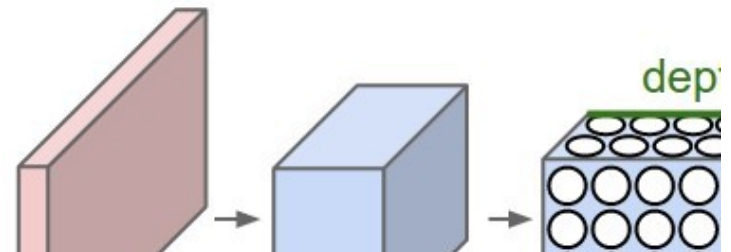
- ❑ Architectures – special graphs with nodes placed on layers
 - Layers
 - ❑ Input layer – size = input's size (#features)
 - ❑ Hidden layers – various sizes (#layers, # neurons/layer)
 - ❑ Output layers – size = output size (e.g. # classes)
 - Topology
 - ❑ Full connected layers (one-way connections, recurrent connections)
- ❑ Mechanism
 - Neuron activation
 - ❑ Constant, step, linear, sigmoid
 - Cost & Loss function → smooth cost function (depends on w&b)
 - ❑ Difference between desired (D) and computed © output
 - ❑ Quadratic cost (*mean squared error*)
 - $\sum ||D - C||^2 / 2n$
 - ❑ Cross-entropy
 - $-\sum [D \ln C + (1 - D) \ln(1 - C)] / n$
 - Learning algorithm
 - ❑ Perceptron rule
 - ❑ Delta rule (Simple/Stochastic Gradient Descent)

Convolutional Neural Networks

- More layers
- More nodes/layer
- Topology of connections
 - Regular NNs → fully connected
 - $O(\text{\#inputs} \times \text{\#outputs})$
 - Conv NNs → partially connected
 - connect each neuron to only a local region of the input volume
 - $O(\text{\#someInputs} \times \text{\#outputs})$



- Topology of layers
 - Regular NNs → linear layers
 - Conv NNs → 2D/3D layers (width, height, depth)



Conv NNs

- Layers of a Conv NN
 - Convolutional Layer → feature map
 - Convolution
 - Activation (thresholding)
 - Pooling/Aggregation Layer → size reduction
 - Fully-Connected Layer → answer

Conv NNs

□ Convolutional layer

■ Aim

- *learn data-specific kernels*
- Perform a linear operation

■ Filters or Local receptive fields or Kernels

□ Content

- Convolution (signal theory) vs. Cross-correlation
- a little (square/cube) window on the input pixels

□ How it works?

- slide the local receptive field across the entire input image

□ Size

- Size of field/filter (F)
- Stride (S)

□ Learning process

- each hidden neuron has
 - FxF shared weights connected to its local receptive field
 - a shared bias
 - an activation function
- each connection learns a weight
- the hidden neuron learns an overall bias as well
- all the neurons in the first hidden layer detect exactly the same feature (just at different locations in the input image) → map from input to the first hidden layer = feature map / activation map

Conv NNs

■ Convolutional Layer – How does it work?

■ Take an input I (example, instance, data) of various dimensions

- A signal \rightarrow 1D input (I_{length})
- a grayscale image \rightarrow 2D input (I_{width} & I_{height})
- an RGB image \rightarrow 3D input (I_{width} , I_{height} & $I_{\text{depth}} = 3$)

■ Consider a set of filters (kernels) $F_1, F_2, \dots, F_{\# \text{filters}}$

- A filter must have the same # dimensions as the input

■ A signal \rightarrow 1D filter

- $F_{\text{length}} \ll I_{\text{length}}$

■ a grayscale image \rightarrow 2D filter

- $F_{\text{width}} \ll I_{\text{width}}$ & $F_{\text{height}} \ll I_{\text{height}}$

■ an RGB image \rightarrow 3D filter

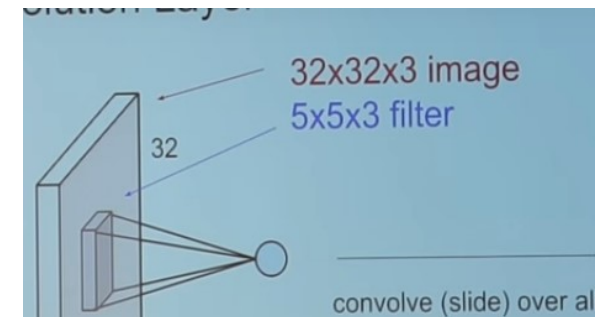
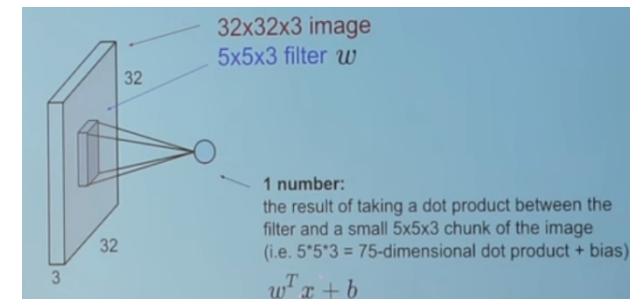
- $F_{\text{width}} \ll I_{\text{width}}$ & $F_{\text{height}} \ll I_{\text{height}}$ &
- $F_{\text{depth}} = I_{\text{depth}} = 3$

■ Apply each filter over the input

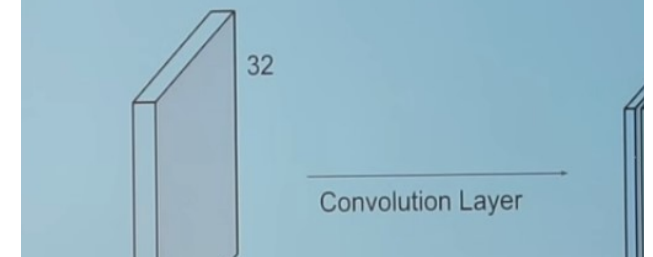
- Overlap filter over a window of the input
 - Stride
 - Padding
- Multiply the filter and the window
- Store the results in an activation map
 - # activation maps = # filters

■ Activate all the elements of each activation map

- ReLU or other activation function



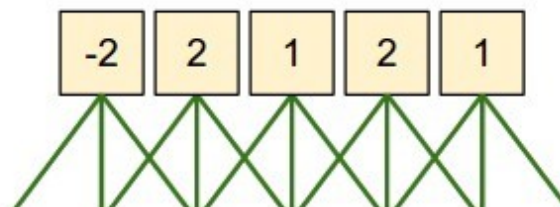
For example, if we had 6 5x5 filters, we'll get 6 se



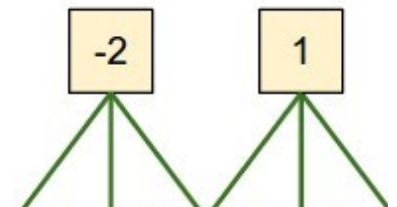
Conv NNs

□ Convolutional layer - Hyperparameters

- input volume size N (L or W_I & H_I or W_I & H_I & D_I)
- size of zero-padding of input volume P (P_L or P_W & P_H or P_W & P_H & P_D)
- the receptive field size (filter size) F (F_L , F_W & F_H , F_W & F_H & F_D)
- stride of the convolutional layer S (S_L , S_W & S_H , S_W & S_H & S_D)
- # of filters (K)
 - depth of the output volume
- # neurons of an activation map = $(N + 2P - F)/S + 1$
- Output size (O or W_O & H_O or W_O & H_O & D_O)
 - $K * [(N + 2P - F)/S + 1]$
- $N = L = 5$, $P = 1$,
 - $F = 3$, $S = 1$
 - $F = 3$, $S = 2$



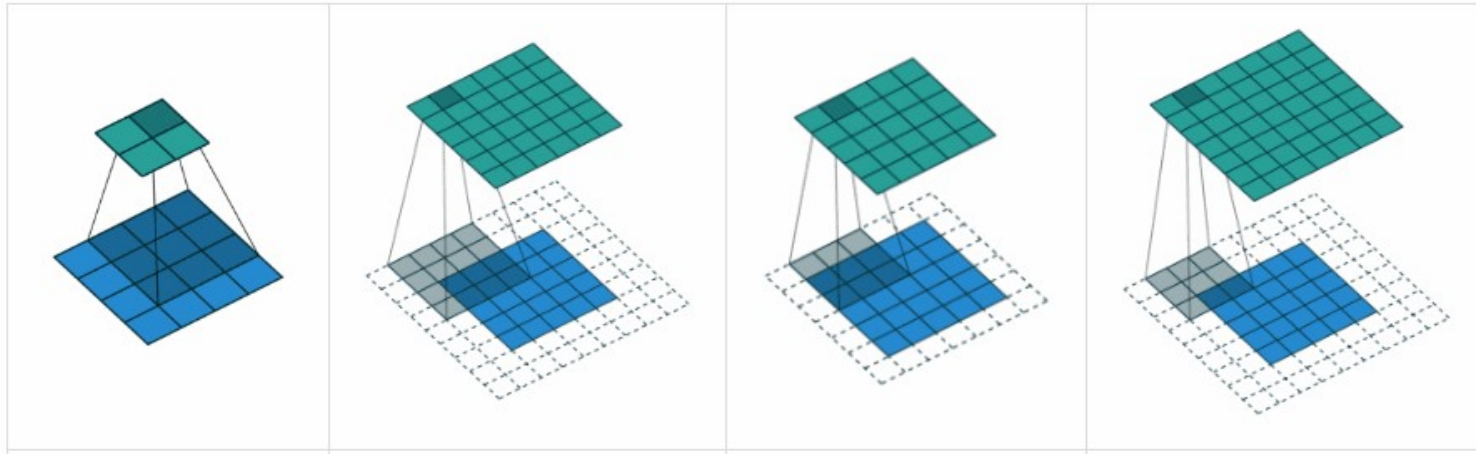
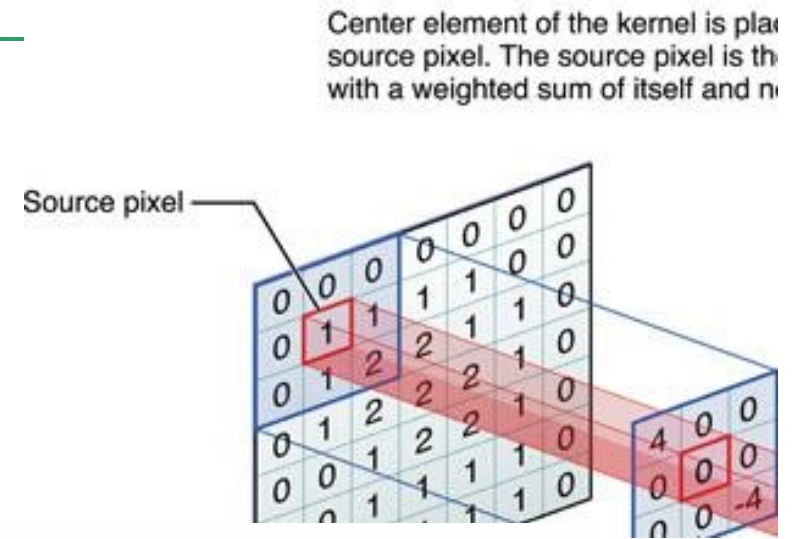
MIRPR - ANN



56

Conv NNs

□ Convolutional Layer – How does it work?



$$\begin{aligned} N &= W_I = H_I = 4, \\ F &= F_W = F_H = 3, \\ P &= 0, S = 1, \\ W_O &= H_O = 2 \end{aligned}$$

$$\begin{aligned} N &= W_I = H_I = 5, \\ F &= F_W = F_H = 4, \\ P &= 2, S = 1, \\ W_O &= H_O = 6 \end{aligned}$$

$$\begin{aligned} N &= W_I = H_I = 5, \\ F &= F_W = F_H = 4, \\ P &= 1, S = 1, \\ W_O &= H_O = 5 \end{aligned}$$

$$\begin{aligned} N &= W_I = H_I = 5, \\ F &= F_W = F_H = 3, \\ P &= 2, S = 1, \\ W_O &= H_O = 7 \end{aligned}$$

Conv NNs

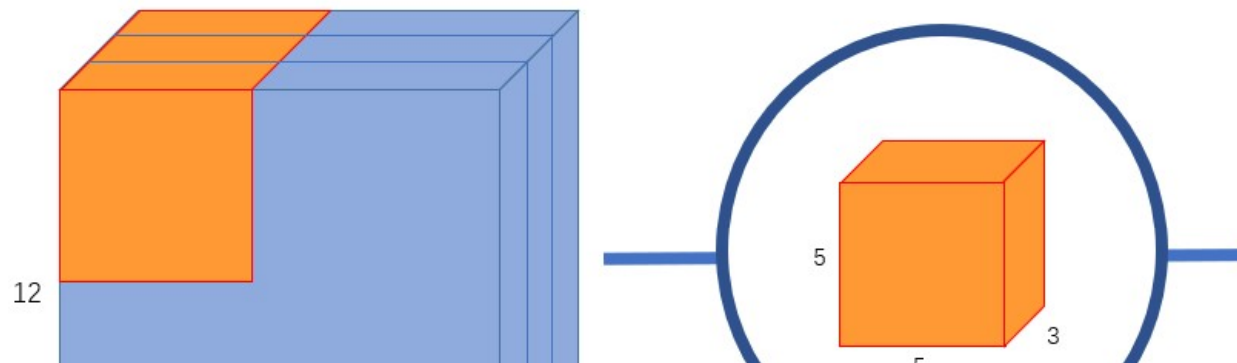
- Convolutional layer – typology
 - Classic convolution
 - Transposed convolution (deconvolution)
 - Dilated convolution
 - Spatial separable (depthwise separable) convolution
 - Grouped convolutions

Conv NNs

□ Convolutional layer – typology

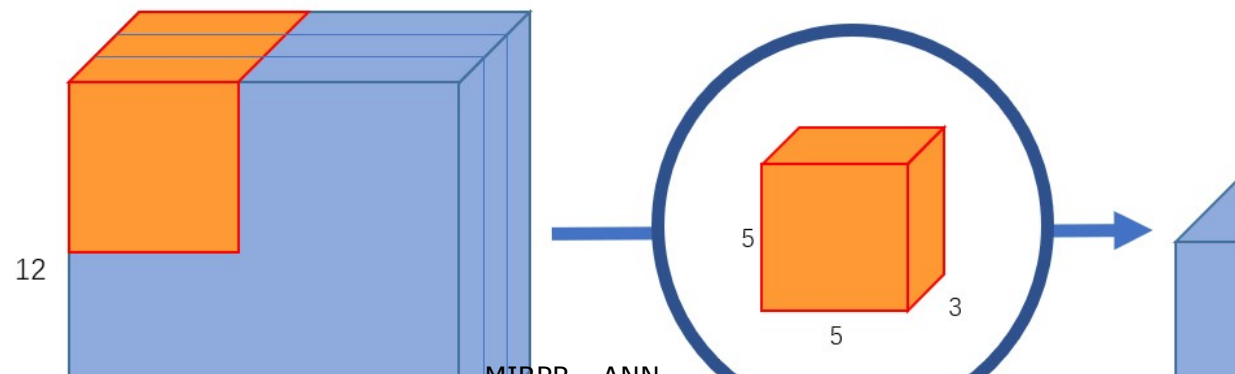
■ classic convolution

- one filter, D channels



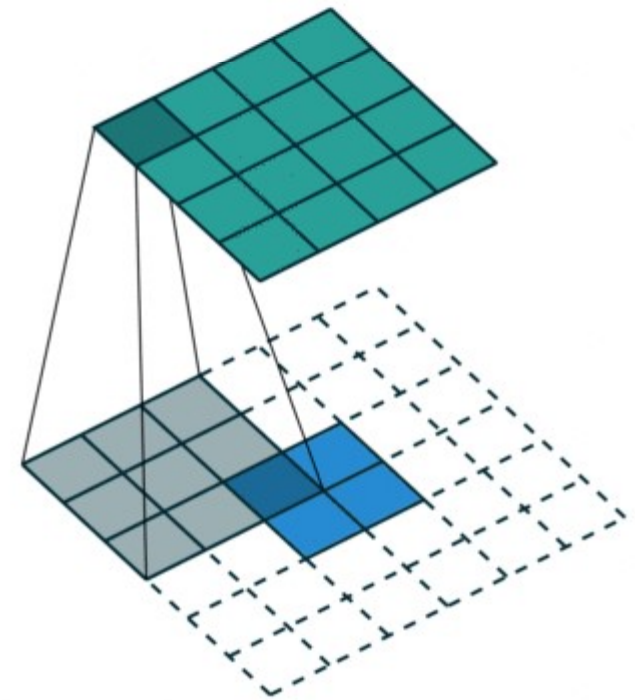
- more filters (K), D channels

* K = * 256



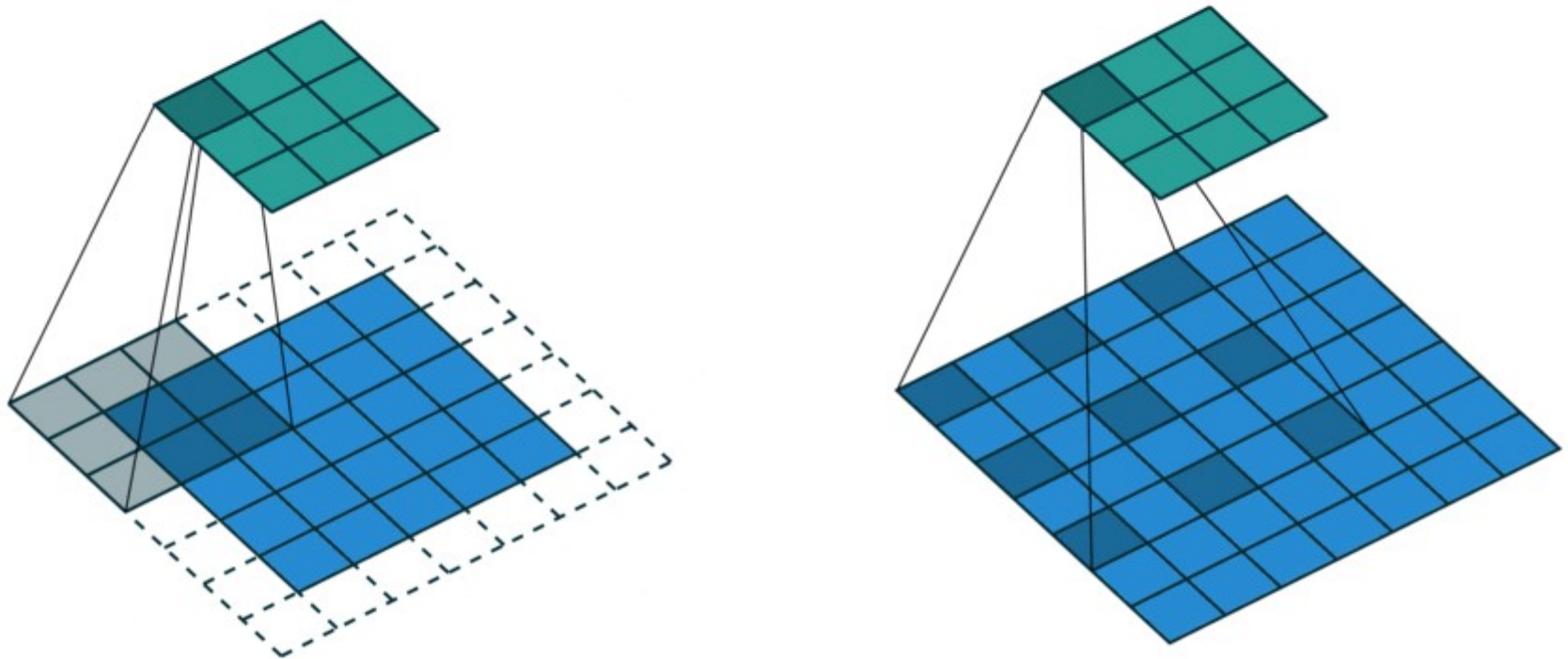
Conv NNs

- Convolutional layer – typology
 - Transposed convolution (deconvolution)
 - See <https://arxiv.org/pdf/1603.07285.pdf>
 - Up-sampling - How?
 - $\text{Img}_{\text{Large}} * F = \text{Img}_{\text{Small}}$
 - $\text{Img}_{\text{Large}} * F * F^T = \text{Img}_{\text{Small}} * F^T$
 - $\text{Img}_{\text{Large}} * I = \text{Img}_{\text{Small}} * F^T$
 - $\text{Img}_{\text{Large}} = \text{Img}_{\text{Small}} * F^T$



Conv NNs

- Convolutional layer – typology
 - Dilated convolutions (atrous convolutions)
 - See <https://arxiv.org/pdf/1511.07122.pdf>



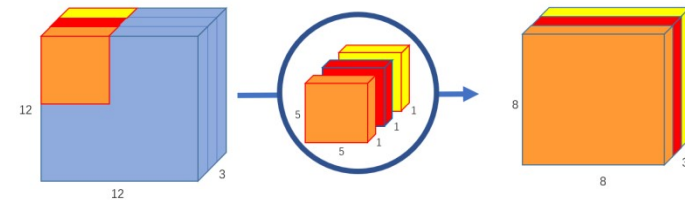
Conv NNs

□ Convolutional layer – typology

■ Spatial separable convolution (Depth-wise separable convolution)

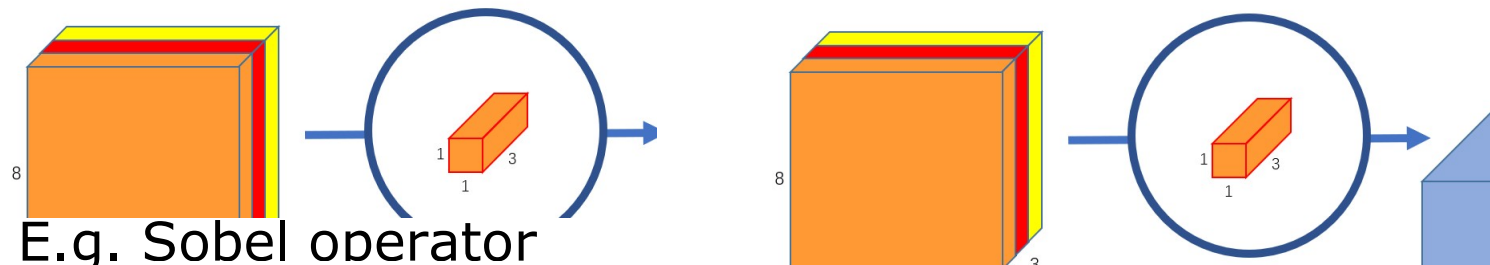
□ Split the convolution into

■ A depthwise convolution



■ A pointwise convolution

* K = * 256



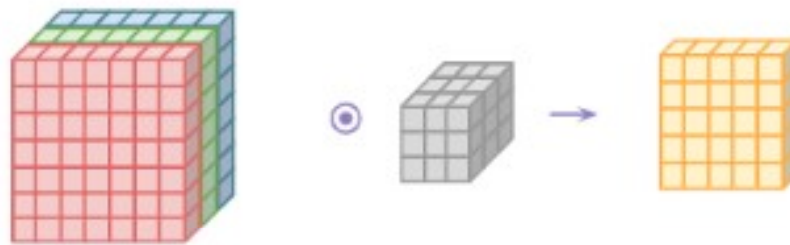
□ E.g. Sobel operator

$$\begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \times$$

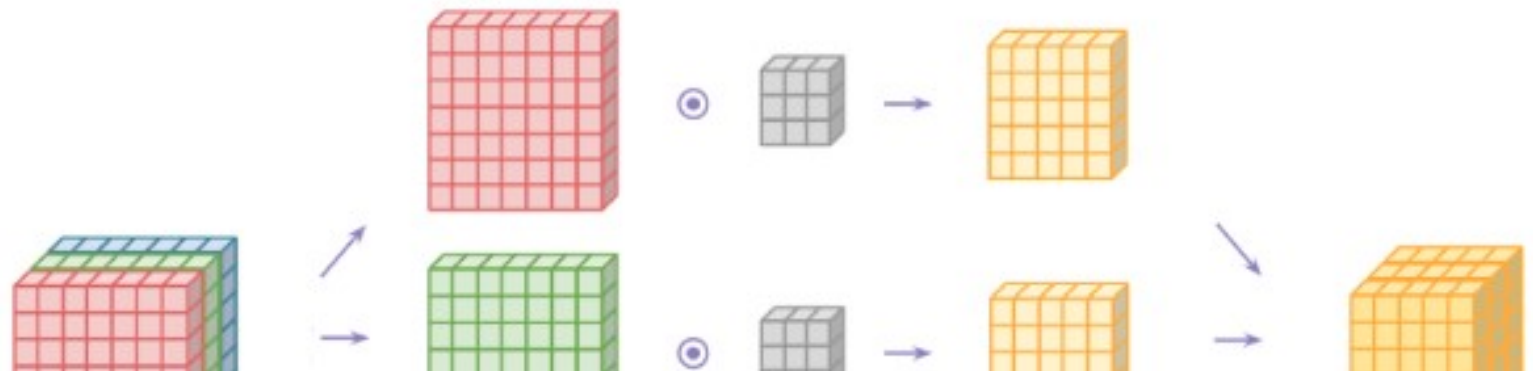
Conv NNs

□ Convolutional layer – typology

- Standard Convolution



- Depthwise Separable Convolution



Conv NNs

□ Convolutional layer – typology

■ Classic vs. spatial separable

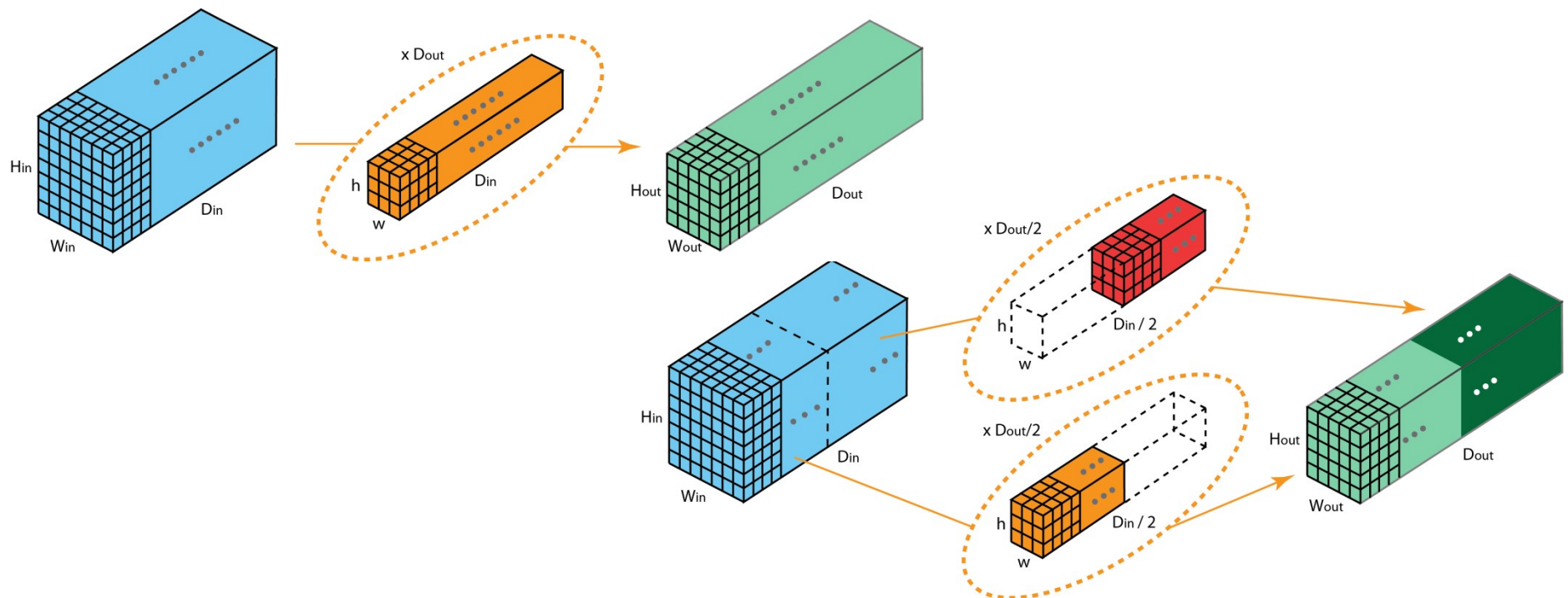
| classic | Spatial separable |
|---|--|
| <p>Image: $W_I \times H_I \times D_I = 12 \times 12 \times 3$</p> <p>1 Filter: $F_W \times F_H \times F_D = 5 \times 5 \times 3$</p> <p>No Padding: $P = 0$ Stride 1: $S = 1$ Output: $W_O \times H_O \times D_O = 8 \times 8 \times 1$</p> <p>$1 * (5 * 5 * 3) * (8 * 8) = 4800$ => 4 800 ops.</p> | <p>Image: $W_I \times H_I \times D_I = 12 \times 12 \times 3$</p> <p>3 Filters: $F_W \times F_H \times F_D = 5 \times 5 \times 1$ 1 Filter: $F'_W \times F'_H \times F'_D = 1 \times 1 \times 3$</p> <p>No Padding: $P = 0$ Stride 1: $S = 1$ Output: $W_O \times H_O \times D_O = 8 \times 8 \times 256$</p> <p>$3 * (5 * 5) * (8 * 8) = 4800$ $1 * (1 * 1 * 3) * (8 * 8) = 192$ => 4 992 ops</p> |
| <p>Image: $W_I \times H_I \times D_I = 12 \times 12 \times 3$</p> <p>256 Filters: $F_W \times F_H \times F_D = 5 \times 5 \times 3$</p> <p>No Padding: $P = 0$ Stride 1: $S = 1$ Output: $W_O \times H_O \times D_O = 8 \times 8 \times 256$</p> <p>$256 * (5 * 5 * 3) * (8 * 8) = 1228800$ => 1 228 800 ops</p> | <p>Image: $W_I \times H_I \times D_I = 12 \times 12 \times 3$</p> <p>3 Filters: $F_W \times F_H \times F_D = 5 \times 5 \times 1$ 256 Filters: $F'_W \times F'_H \times F'_D = 1 \times 1 \times 3$</p> <p>No Padding: $P = 0$, Stride 1: $S = 1$ Output: $W_O \times H_O \times D_O = 8 \times 8 \times 256$</p> <p>$3 * (5 * 5) * (8 * 8) = 4800$ $256 * (1 * 1 * 3) * (8 * 8) = 49152$ => 53 952 ops</p> |

Conv NNs

□ Convolutional layer – typology

■ Grouped convolutions

- See <https://arxiv.org/pdf/1605.06489.pdf>
- Efficient training (More GPUs) => model parallelisation
- Fewer parameters
- Better representations



Conv NNs

- Convolutional layer – typology
 - Classic vs. grouped convolutions

| Classic | Group |
|--|---|
| <p>Image: $W_I \times H_I \times D_I = 12 \times 12 \times D_I$</p> <p>$D_O$ Filter: $F_W \times F_H \times D_I = 5 \times 5 \times D_I$</p> <p>No Padding: $P = 0$ Stride 1: $S = 1$ Output: $W_O \times H_O \times D_O = 8 \times 8 \times D_O$</p> <p>#operations: $F_W \times F_H \times D_I \times D_O = 5 \times 5 \times D_I \times D_O$</p> | <p>Image: $W_I \times H_I \times D_I = 12 \times 12 \times D_I$</p> <p>$D_O/2$ Filters: $F_W \times F_H \times (D_I / 2) = 5 \times 5 \times (D_I / 2)$ $D_O/2$ Filters: $F_W \times F_H \times (D_I / 2) = 5 \times 5 \times (D_I / 2)$</p> <p>No Padding: $P = 0$ Stride 1: $S = 1$ Output: $W_O \times H_O \times D_O = 8 \times 8 \times D_O$</p> <p>#operations: $[F_W \times F_H \times (D_I / 2) \times (D_O / 2)] * 2$ $= 5 \times 5 \times D_I \times D_O / 2$</p> |

Conv NNs

- Convolutional layer - ImageNet challenge in 2012 (Alex Krizhevsky <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>)
 - Input images of size $[227 \times 227 \times 3]$
 - $F=11, S=4, P=0, K=96 \rightarrow$ Conv layer output volume of size $[55 \times 55 \times 96]$
 - $55 \times 55 \times 96 = 290,400$ neurons in the first Conv Layer
 - each has $11 \times 11 \times 3 = 363$ weights and 1 bias.
 - $290400 * 364 = 105,705,600$ parameters on the first layer

Conv NNs

□ Convolutional layer - parameter sharing scheme

- constrain the neurons in each depth slice to use the same weights and bias
 - detect exactly the same feature, just at different locations in the input image
 - convolutional networks are well adapted to the translation invariance of images
- Example
 - 96 unique set of weights (one for each depth slice),
 - for a total of $96 \cdot 11 \cdot 11 \cdot 3 = 34,848$ unique weights,
 - 34,944 parameters (+96 biases).
- if all neurons in a single depth slice are using the same weight vector, then the forward pass of the CONV layer can in each depth slice be computed as a **convolution** of the neuron's weights with the input volume → the name: Convolutional Layer
- Set of weights = filter (kernel)
- Can be applied, but are image-dependent
 - faces that have been centered in the image
- A Convolutional Layer without parameter sharing → **Locally-Connected Layer**

Convolutional Neural Networks

- CNN important strong points:
 - Sparse interactions
 - Parameter sharing
 - Equivariant representations
 - Pro: Translation
 - Cons: rotation, scale
 - allows for working with inputs of variable size

Conv NNs

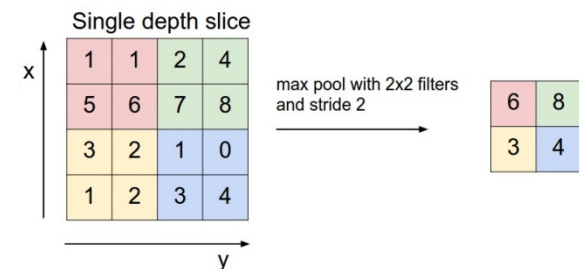
□ Pooling layer

■ Aim

- progressively reduce the spatial size of the representation
 - to reduce the amount of parameters and computation in the network
 - to also control overfitting
- a subsampling step
 - downsample the spatial dimensions of the input.
- simplify the information in the output from the convolutional layer

■ How it works

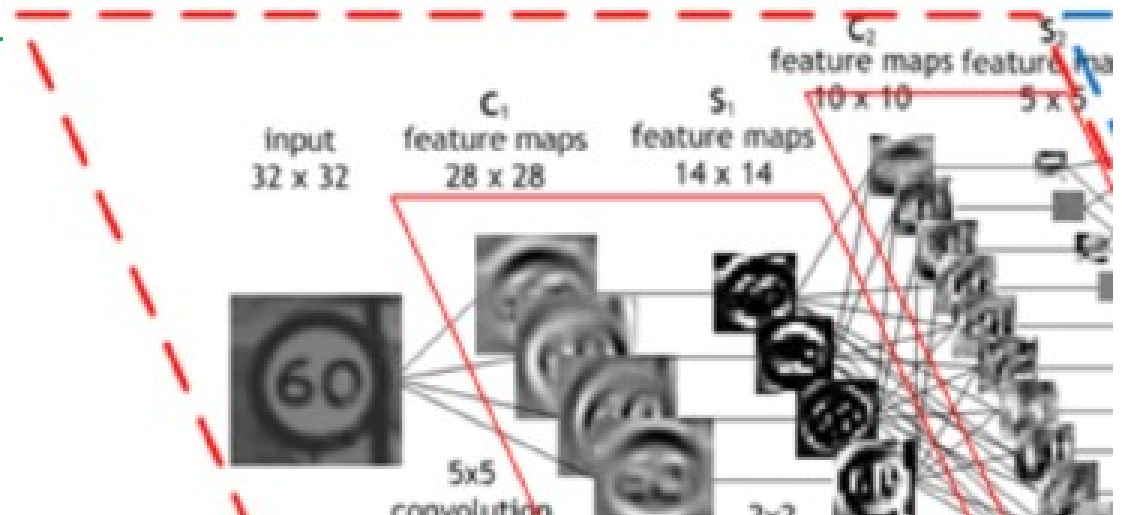
- takes each feature map output from the convolutional layer and prepares a condensed feature map
- each unit in the pooling layer may summarize a region in the previous layer
- apply pooling filters to each feature map separately
 - Pooling filter size (spatial extent of pooling) PF
 - Pooling filter stride PS
 - No padding



Conv NNs

□ Pooling layer

- How it works



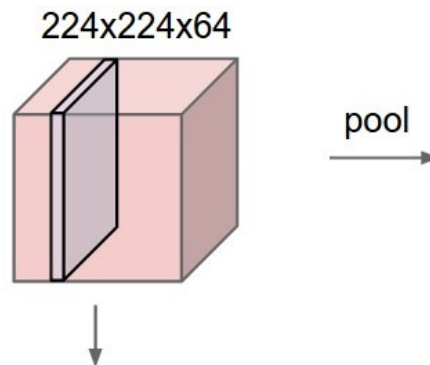
- resizes it spatially, using
 - the MAX operation
 - the average operation
 - L^p norm: $\sqrt[p]{\sum x^p}$
 - L²-norm operation (square root of the sum of the squares of the activations in a rectangular neighbourhood/region) $\leftrightarrow p = 2$
- Log prob PROB: $\frac{1}{b} \log(\sum e^{bx})$

Conv NNs

□ Pooling layer

■ Two reasons:

□ Dimensionality reduction



□ Invariance to transformation (rotation, translation)

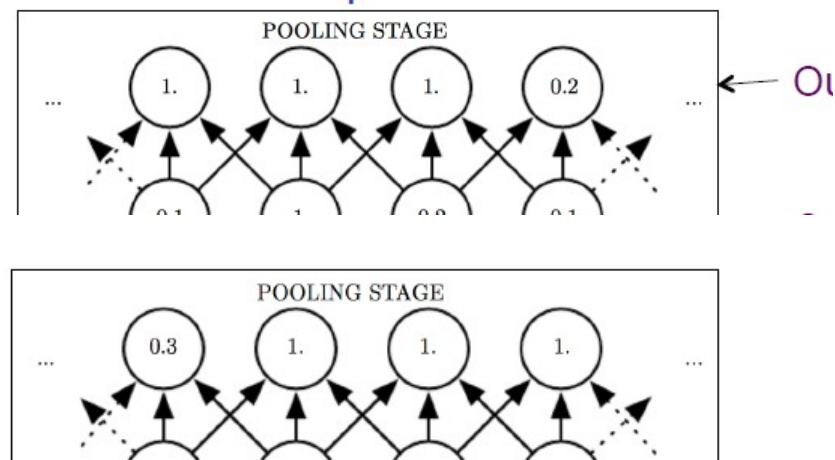


Conv NNs

□ Pooling layer

■ Two reasons:

- Invariance to transformation (rotation, translation)
 - Small translations – e.g. Max pooling



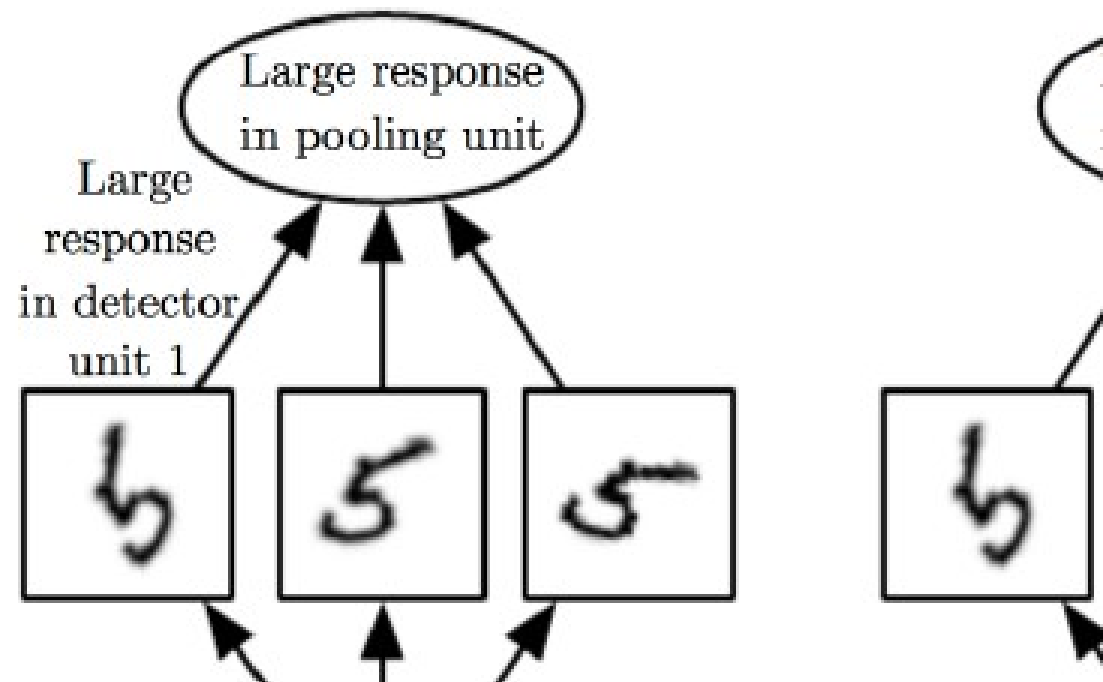
- When? => if we care about whether a feature is present rather than exactly where it is

Conv NNs

□ Pooling layer

■ Two reasons:

- Invariance to transformation (rotation, translation)
 - Rotations



Conv NNs

□ Pooling layer

■ Size conversion

- Input:
 - $K \times N$
- Output
 - $K \times [(N - PF) / PS + 1]$

■ Typology

- Local pooling (patch-based pooling)
- Global pooling (image-based pooling)

■ Remark

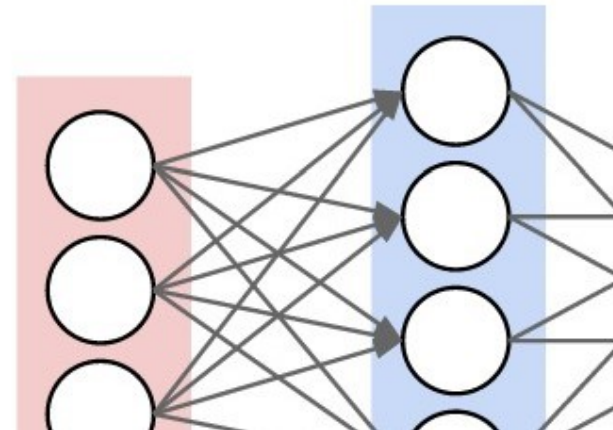
- introduces zero parameters since it computes a fixed function of the input
- note that it is not common to use zero-padding for Pooling layers
- pooling layer with $PF=3, PS=2$ (also called overlapping pooling), and more commonly $PF=2, PS=2$
- pooling sizes with larger filters are too destructive
- keep track of the index of the max activation (sometimes also called *the switches*) so that gradient routing is efficient during backpropagation

Conv NNs

□ Fully-connected layer

- Neurons have full connections to all inputs from the previous layer

- Various activations
 - ReLU (often)



Conv NNs

□ CNN architectures

- INPUT \rightarrow $[[\text{CONV} \rightarrow \text{RELU}]^N \rightarrow \text{POOL?}]^M \rightarrow [\text{FC} \rightarrow \text{RELU}]^K \rightarrow \text{FC}$
- Most common:
 - INPUT \rightarrow FC \leftrightarrow a linear classifier
 - INPUT \rightarrow CONV \rightarrow RELU \rightarrow FC
 - INPUT \rightarrow $[\text{CONV} \rightarrow \text{RELU} \rightarrow \text{POOL}]^2 \rightarrow \text{FC} \rightarrow \text{RELU} \rightarrow \text{FC}$.
 - INPUT \rightarrow $[\text{CONV} \rightarrow \text{RELU} \rightarrow \text{CONV} \rightarrow \text{RELU} \rightarrow \text{POOL}]^3 \rightarrow [\text{FC} \rightarrow \text{RELU}]^2 \rightarrow \text{FC}$
 - a good idea for larger and deeper networks, because multiple stacked CONV layers can develop more complex features of the input volume before the destructive pooling operation

Conv NNs

□ Remarks

- Prefer a stack of small filter CONV to one large receptive field CONV layer
 - Pro:
 - Non-linear functions
 - Few parameters
 - Cons:
 - more memory to hold all the intermediate CONV layer results
- Input layer size \rightarrow divisible by 2 many times
- Conv layers \rightarrow small filters
 - $S \geq 1$
 - $P = (F - 1) / 2$
- Pool layers
 - $F \leq 3, S = 2$

Conv NNs

□ Output layer

■ Multiclass SVM

- Largest score indicates the correct answer

■ Softmax (normalized exponential function)

- Largest probability indicates the correct answer
- converts raw scores to probabilities
- "squashes" a $\#classes$ -dimensional vector \mathbf{z} of arbitrary real values to a $\#classes$ -dimensional vector $\sigma(\mathbf{z})$ of real values in the range (0, 1) that add up to 1
- $\sigma(\mathbf{z})_j = \exp(z_j) / \sum_{k=1..\#classes} \exp(z_k)$

Conv NNs

□ Output layer → Multiclass SVM

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_j - s_{y_i} \leq 0 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

“Hinge

| | | | |
|------|------------|------------|-------------|
| cat | 3.2 | 1.3 | 2.2 |
| car | 5.1 | 4.9 | 2.5 |
| frog | -1.7 | 2.0 | -3.1 |



| | | |
|-----|------------|------------|
| cat | 3.2 | 1.3 |
| car | 5.1 | 4.9 |

Conv NNs

□ Output layer → Multiclass SVM

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_j \leq s_{y_i} \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

“Hinge

| | | | |
|------|------------|------------|-------------|
| cat | 3.2 | 1.3 | 2.2 |
| car | 5.1 | 4.9 | 2.5 |
| frog | -1.7 | 2.0 | -3.1 |



| | | |
|-----|------------|------------|
| cat | 3.2 | 1.3 |
| car | 5.1 | 4.9 |

Conv NNs

□ Output layer → Multiclass SVM

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_j \leq s_{y_i} \\ s_j - s_{y_i} + 1 & \text{if } s_j > s_{y_i} \end{cases}$$

"Hinge"

| | | | |
|------|------------|------------|-------------|
| cat | 3.2 | 1.3 | 2.2 |
| car | 5.1 | 4.9 | 2.5 |
| frog | -1.7 | 2.0 | -3.1 |

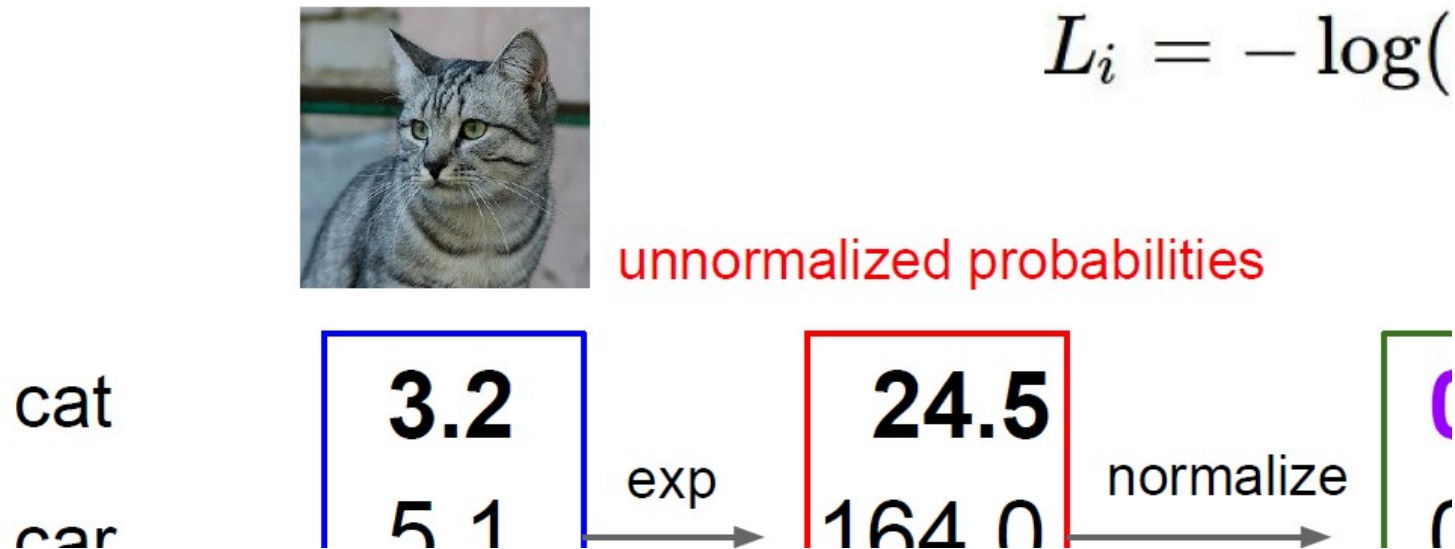


| | | |
|-----|------------|------------|
| cat | 3.2 | 1.3 |
| car | 5.1 | 4.9 |

$$L = (2.9 + 0 + \dots)$$

Conv NNs

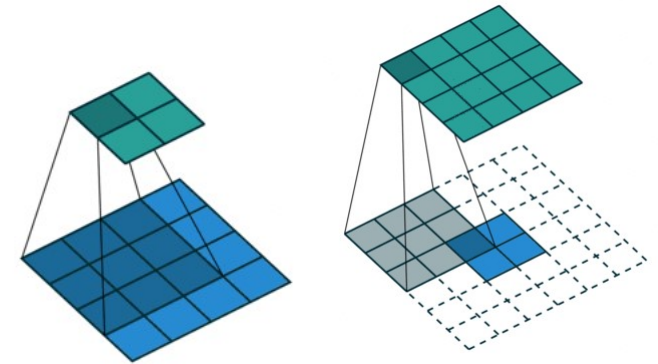
□ Output layer → Softmax



Conv NNs

□ Visualising a CNN

- Convolution vs. deconvolution
- Activation vs. rectification
- Pooling vs. unpooling



| | | | |
|------|------|------|------|
| 0.1 | 0.5 | 1.2 | -0.7 |
| 0.8 | -0.2 | -0.5 | 0.3 |
| 0.4 | 0.9 | -0.1 | -0.2 |
| -0.6 | 0.1 | 0.5 | 0.3 |

max-pooling

| | |
|-----|-----|
| 0.8 | 1.2 |
| 0.9 | 0.5 |

| | | | |
|---|---|-----|---|
| 0 | 0 | 0.5 | 0 |
|---|---|-----|---|

Image classification

□ Task

Classification: C classes

Input: Image

Output: Class label

Evaluation metric: Accuracy



□ Databases

- Pascal VOC <http://host.robots.ox.ac.uk/pascal/VOC/>
 - 2005 – image classification task (4 classes, 1578 images, 2209 objects)
 - 2006 – image classification task (10 classes, 2618 images, 4754 objects)
 - ...
 - 2012 – image classification task (20 classes, 11 530 images, 6929 objects)
- ImageNet <http://www.image-net.org/>
 - 2010 – image classification task only (1000 classes, 14,197,122 images,)
 - 2011, ... - other tasks (localisation, segmentation, detection)

□ Deep Learning Algorithms

- Various CNN architectures (LeNet, AlexNet, VGG, Inception, ResNet, ...)

Conv NNs

□ Common architectures

- LeNet (Yann LeCun, 1998) - <http://yann.lecun.com/exdb/publis/pdf/lecun-98.pdf>
 - A conv layer + a pool layer
- AlexNet (Alex Krizhevsky, Ilya Sutskever and Geoff Hinton, 2012) <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>
 - More conv layers + more pool layers
- ZF Net (Matthew Zeiler and Rob Fergus, 2013) <https://arxiv.org/pdf/1311.2901.pdf>
 - AlexNet + optimisation of hyper-parameters
- GoogleLeNet (Christian Szegedy et al., 2014) <https://arxiv.org/pdf/1409.4842.pdf>
 - *Inception Module* that dramatically reduced the number of parameters in the network (AlexNet 60M, GoogleLeNet 4M) <https://arxiv.org/pdf/1602.07261.pdf>
 - uses Average Pooling instead of Fully Connected layers at the top of the ConvNet → eliminating parameters
- VGGNet (Karen Simonyan and Andrew Zisserman, 2014) <https://arxiv.org/pdf/1409.1556.pdf>
 - 16 Conv/FC layers (FC → a lot more memory; they can be eliminated)
 - pretrained model is available for plug and play use in Caffe
- ResNet (Kaiming He et al., 2015) <https://arxiv.org/pdf/1512.03385.pdf> (Torch)
 - *skip connections*
 - batch normalization

Conv NNs

- Classical CNNs
 - LeNet (1998)
 - AlexNet (2012) – first Deep CNN
 - ZfNet (2013)
- Modern CNNs
 - VGG (2014)
 - NiN (2014)
 - GoogleLeNet (2014)
 - MobileNet (2017)
 - ResNet (2015)

Conv CNNs

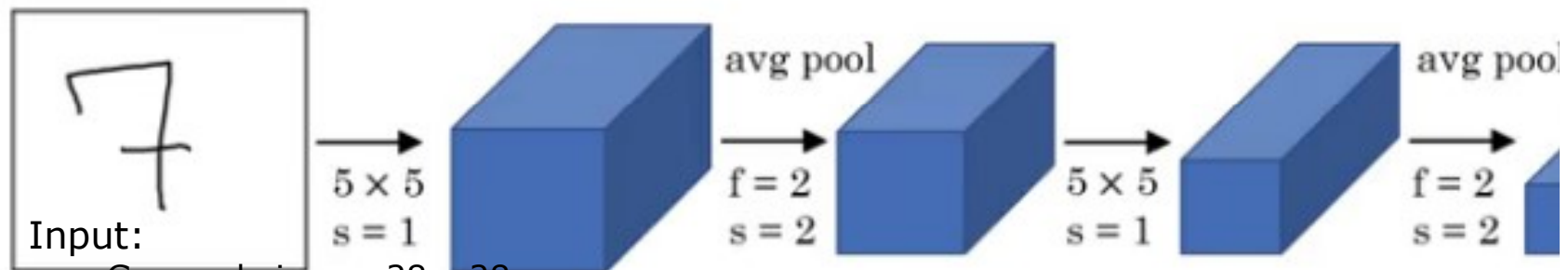
□ Classical CNNs → LeNet (1998)

■ Parent

- Yann LeCun (NYU), MNIST data
- LeCun, Yann, Léon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. "Gradient-Based Learning Applied to Document Recognition." In *Proceedings of the IEEE*, 2278–2324 http://vision.stanford.edu/cs598_spring07/papers/Lecun98.pdf

■ Flow:

- Input → 2 x [Conv → Pool] → FC → FC → softmax → Output (10 classes)



■ Input:

- Grayscale Image 28 x 28

■ Activation:

- Tahn, sigmoid

■ Filters (#filters(size, padding, stride))

- 6(5 x 5, 2, 1), 16(5 x 5, 0, 1)

■ Pooling

- Avg-pooling 2x2, stride 2

■ Loss

- Softmax (cross-entropy)

■ # parameters

- 60 000

Conv CNNs

□ Classical CNNs → LeNet (1998)

■ Flow:

□ Input → 2 x [Conv → Pool] → FC → FC → softmax → Output

■ Activation:

□ Tanh

■ 0 centered → on avg., values → 0 → derivative → 1

□ Sigm

■ 0.5 centered → on avg., values → 0.5 → derivative → $\sim 0.25 < 1$

□ Issues :

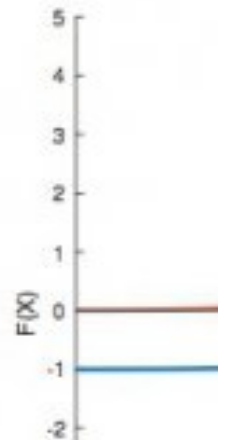
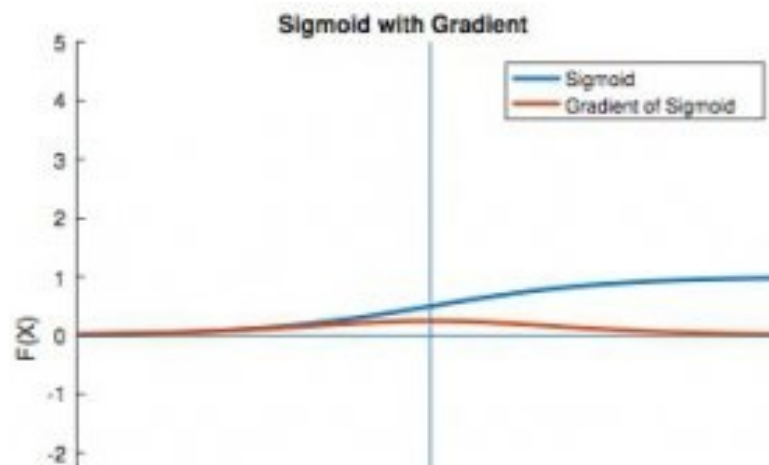
■ vanishing gradient problem (VGP = the gradient of the activation becomes negligible)

$$\text{sigm}(x) = \frac{1}{1+e^{-x}} \in (0,1)$$

$$\text{sigm}'(x) = \text{sigm}(x)(1 - \text{sigm}(x)) \in (0,0.5)$$

$$\text{tahn}(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \in (-1,1)$$

$$\text{tahn}'(x) = 1 - (\text{tahn}(x))^2 \in (0,1)$$



Conv NNs

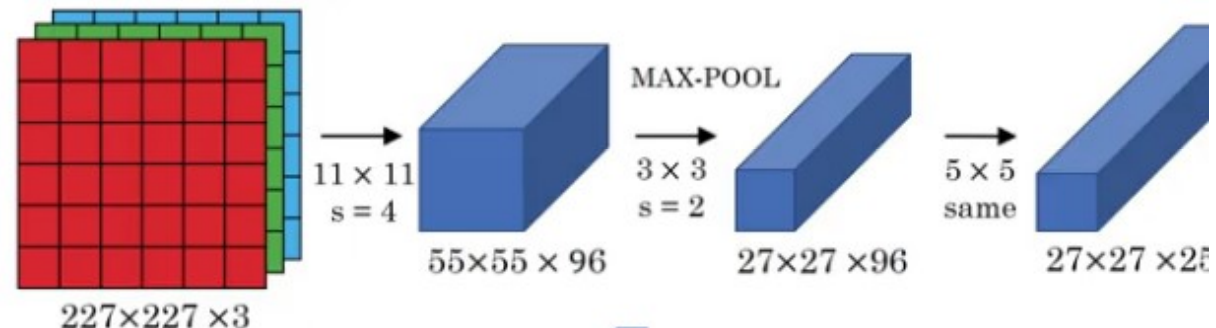
□ Classic CNNs → AlexNet (2012)

■ Parents

- Alex Krizhevsky et al., ImageNet data
- Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. 2012. "ImageNet Classification with Deep Convolutional Neural Networks." In *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1*, 1097–1105. NIPS'12. Lake Tahoe, Nevada <https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>
- Winner of ILSVRC 2012 (5-6 days for training - GTX 580 GPUs)

■ Flow:

- Input → 2 x [Conv → Pool → **Norm**] → 3 x Conv → Pool → 3 x [FC → **DropOut**] → Softmax → Output (1000 classes)



■ Input

- RGB image 224 x 224 x 3

■ Activation

- ReLU, sigmoid

■ Filters (#filters(size, padding, stride))

- 96(11 x 11, 0, 4), 256(5 x 5, 2, 1), 3 x [384(3 x 3, 1, 1)]

■ Pooling

- Max-pooling 3x3, stride 2

■ Loss

- Cross-entropy loss

■ # parameters

- 60 000 000

Conv NNs

□ Classic CNNs → AlexNet (2012)

■ Flow:

□ Input → 2 x [Conv → Pool → **Norm**] → 3 x Conv → Pool → 3 x [FC → **DropOut**] → Softmax → Output

■ Activation

□ Conv → ReLU

■ Advantages

- Feature sparsity
- Reducing VGP

■ Drawbacks

- Dying ReLU problem: $\text{output}(\text{node}) < 0 \rightarrow \text{derivative} = 0 \rightarrow \text{weights are not changed / trained}$

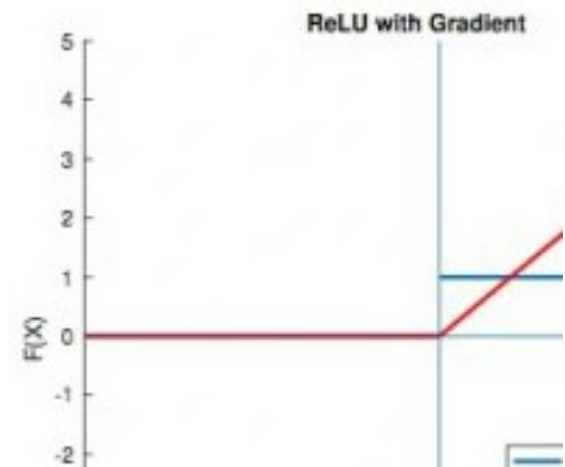
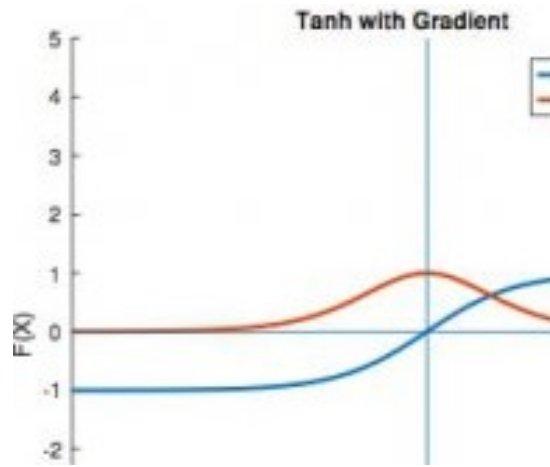
- Conv layers are more affected by VGP

□ FC → tanh

- FC layers are less affected by VGP

$$\text{ReLU}(x) = \begin{cases} 0, & \text{for } x \leq 0 \\ x, & \text{for } x > 0 \end{cases}$$

$$\text{ReLU}'(x) = \begin{cases} 0, & \text{for } x \leq 0 \\ 1, & \text{for } x > 0 \end{cases}$$



Conv NNs

□ Classic CNNs → AlexNet (2012)

■ Flow:

- Input → 2 x [Conv → Pool → **Norm**] → 3 x Conv → Pool → 3 x [FC → **DropOut**] → Softmax → Output

■ Normalisation layers

- normalize the activations of each node by subtracting its mean and dividing by its standard deviation estimating both quantities based on the statistics of the current the current minibatch
- typically applied BN after the convolution and before the nonlinear activation function
- Applied on each channel / feature map

$$BN(x) = \gamma \frac{x - \mu}{\sigma} + \beta, \quad \mu = \frac{1}{|batch|} \sum_{x \in batch} x, \quad \sigma^2 = \frac{1}{|batch|} \sum_{x \in batch} (x - \mu)^2 + \varepsilon$$

■ Dropout layers

- Help in removing complex co-adaptations (reducing overfitting)
 - #training samples > 10 * # parameters
- Net is more robust to noise

Conv NNs

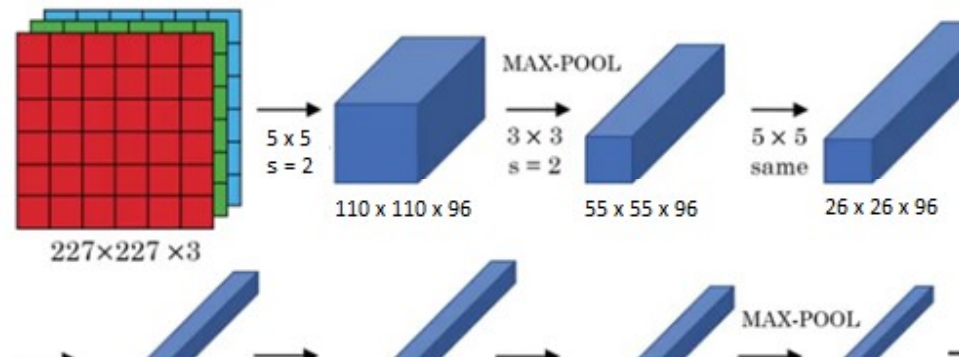
□ Classic CNNs → ZfNet (2013)

■ Parents

- Rob Fergus, Matthew D. Zeiler (NYU) → ClarifAI, CIFAR-10
- Zeiler, Matthew D., and Rob Fergus. 2013. "Visualizing and Understanding Convolutional Networks." *CoRR* abs/1311.2901 <http://arxiv.org/abs/1311.2901>
- Winner of ILSVRC 2013

■ Flow

- AlexNet improved based on visualisation of the feature maps



■ Input

- RGB image 224 x 224 x 3

■ Activation

- ReLU, tahn

■ Filters (#filters(size, padding, stride))

- 96(7 x 7, 0, 2), 256(5 x 5, 2, 1), 512(3 x 3, 1, 1), 1024(3 x 3, 1, 1), 512(3 x 3, 1, 1)

■ Pooling

- Max-pooling 3x3, stride 2

■ Loss

- Cross-entropy

■ # parameters

- TBA

Reducing the filter' size and stride → More mid-frequencies

Increase the # filters / feature maps → more features

Conv NNs

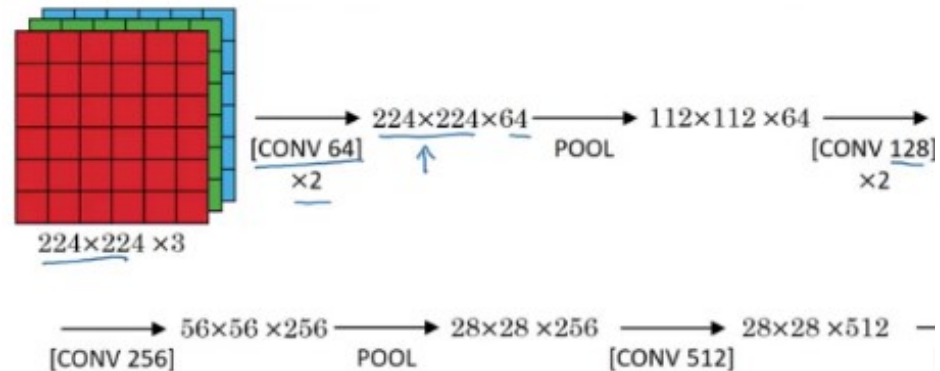
Modern CNNs → VGG (2014)

Parents

- Simonyan and Zisserman (Visual Geometry Group, Oxford)
- Simonyan, Karen, and Andrew Zisserman. 2014. "Very Deep Convolutional Networks for Large-Scale Image Recognition." *CoRR* abs/1409.1556. <http://arxiv.org/abs/1409.1556>.
- Deeper is better!
- Winner of ILSVRC 2014 (localisation), 2nd place (classification)

Flow

- Input → 2 x [Conv → Pool] → 3 x [Conv → Conv → Pool] → FC → FC → FC → Softmax → Output (1000 classes) → VGG-11
- Input → 3 x [Conv → Conv → Pool] → 2 x [Conv → Conv → Conv → Pool] → FC → FC → FC → Softmax → Output → VGG-16
- Input → 3 x [Conv → Conv → Pool] → 2 x [Conv → Conv → Conv → Conv → Pool] → FC → FC → FC → Softmax → Output → VGG-19



Input

- RGB image 224 x 224 x 3

Activation

- ReLU, tahn

Blocks = patterns of layers (Conv + ReLU + Pool)

- Stacked smaller filters → complex features learnt at a lower cost

Loss

- Cross-entropy loss

parameters

- 138 000 000 (VGG16)

Conv NNs

□ Modern CNNs → VGG (2014)

■ Parents

- Simonyan and Zisserman (Visual Geometry Group, Oxford)
- Simonyan, Karen, and Andrew Zisserman. 2014. "Very Deep Convolutional Networks for Large-Scale Image Recognition." *CoRR* abs/1409.1556. <http://arxiv.org/abs/1409.1556>.
- **Deeper is better!**
- Winner of ILSVRC 2014 (localisation), 2nd place (classification)

■ Flow

- Input → 2 x [Conv → Pool] → 3 x [Conv → Conv → Pool] → FC → FC → FC → Softmax → Output (1000 classes) → VGG-11
- Input → 3 x [Conv → Conv → Pool] → 2 x [Conv → Conv → Conv → Pool] → FC → FC → FC → Softmax → Output → VGG-16
- Input → 3 x [Conv → Conv → Pool] → 2 x [Conv → Conv → Conv → Conv → Pool] → FC → FC → FC → Softmax → Output → VGG-19

■ Blocks = patterns of layers (Conv + ReLU + Pool)

- Stacked smaller filters → complex features learnt at a lower cost
- [Conv64(3 x 3, 1, 1) → ReLU → Max-pooling 3x3, stride 2]
- [Conv128(3 x 3, 1, 1) → ReLU → Max-pooling 3x3, stride 2]
- [2 x Conv256(3 x 3, 1, 1) → ReLU → Max-pooling 2x2, stride 2]
- [2 x Conv512(3 x 3, 1, 1) → ReLU → Max-pooling 2x2, stride 2]
- [2 x Conv512(3 x 3, 1, 1) → ReLU → Max-pooling 2x2, stride 2]

■ Filters

- Stacks of smaller filters
 - 3 Conv(3x3,0,1) ↔ 1 Conv(7 x 7, 0, 1)
 - Deeper → more non-linearities
 - Fewer parameters ($3 * 3^2 * \text{\#Channels}^2 \leftrightarrow 1 * 7^2 * \text{\#Channels}^2$)

■ # parameters

- 138 000 000 (VGG16) – a large part of them are used in the final 3 Fully Connected layers
 - First FC layer → AvgPooling (InceptionNet, ResNet)

Conv NNs

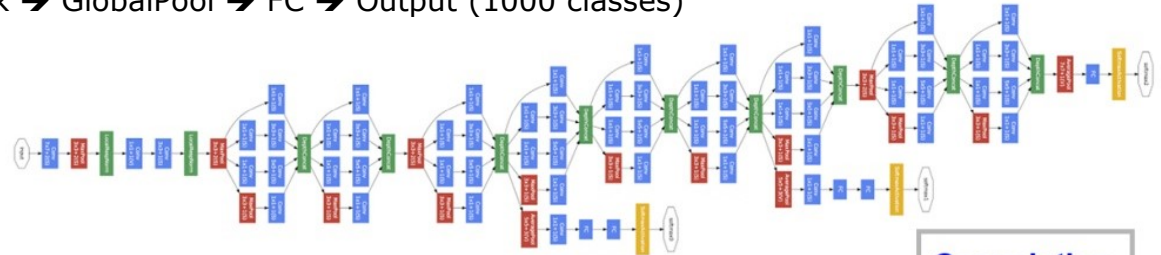
Modern CNNs → GoogleLeNet / InceptionNet(2014)

■ Parents

- Szegedy, Christian, Sergey Ioffe, and Vincent Vanhoucke. 2016. "Inception-V4, Inception-Resnet and the Impact of Residual Connections on Learning." *CoRR* abs/1602.07261. <https://arxiv.org/pdf/1409.4842.pdf>
- Inception-V4 <http://arxiv.org/abs/1602.07261>
- Top performamnce ILSVRC 2014 (winner of classification)
- movie Inception ("We Need To Go Deeper")

■ Flow

- Input → Conv → Pool → Conv → Conv → Pool → 2 x InceptionBlock → Pool → 5 x InceptionBlock → Pool → 2 x InceptionBlock → GlobalPool → FC → Output (1000 classes)



■ Particularities

- FC layers are replaced by AvgPooling → reducing #parameters
- Repeatedly usage of the inception block → multiple filters of different sizes.

■ Input

- RGB image 224 x 224 x 3

■ Activation

- ReLU, than

■ Loss

- Softmax (cross-entropy)

■ # parameters

- 5 000 000 parameters (12x less than AlexNet)

Conv NNs

Modern CNNs → GoogleLeNet / InceptionNet(2014)

Flow

- Input → Conv → Pool → Conv → Conv → Pool → 2 x InceptionBlock → Pool → 5 x InceptionBlock → Pool → 2 x InceptionBlock → GlobalPool → FC → Output (1000 classes)

Particularities

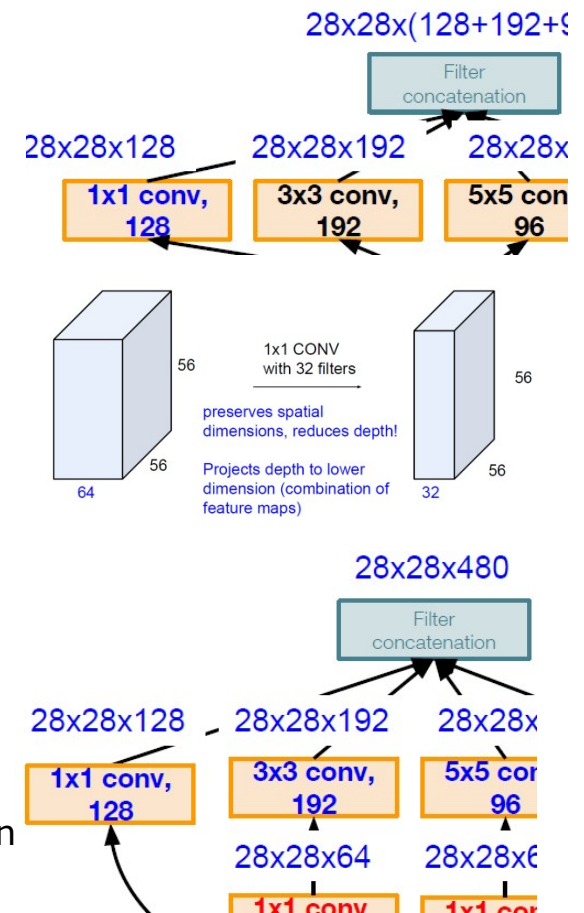
Naïve Inception block/module

- More Conv Ops:
 - [1x1 conv, 128] $28 \times 28 \times 128 \times 1 \times 1 \times 256$
 - [3x3 conv, 192] $28 \times 28 \times 192 \times 3 \times 3 \times 256$
 - [5x5 conv, 96] $28 \times 28 \times 96 \times 5 \times 5 \times 256$
- Total: 854M ops, 592 K param

Reduced Inception module

- Use 1x1 conv to reduce feature depth**
- [1x1 conv, 64] $28 \times 28 \times 64 \times 1 \times 1 \times 256$
- [1x1 conv, 64] $28 \times 28 \times 64 \times 1 \times 1 \times 256$
- [1x1 conv, 128] $28 \times 28 \times 128 \times 1 \times 1 \times 256$
- [3x3 conv, 192] $28 \times 28 \times 192 \times 3 \times 3 \times 64$
- [5x5 conv, 96] $28 \times 28 \times 96 \times 5 \times 5 \times 64$
- [1x1 conv, 64] $28 \times 28 \times 64 \times 1 \times 1 \times 256$
- Total: 358M ops, 376 K param

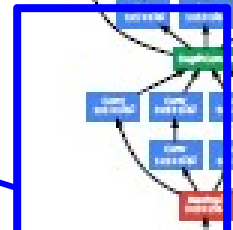
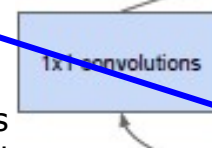
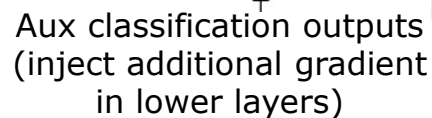
- InceptionNet = Stack of Inception module with dimension



| Age Group | Male (%) | Female (%) | Total (%) |
|-----------|----------|------------|-----------|
| 18-24 | 45 | 55 | 50 |
| 25-34 | 40 | 60 | 50 |
| 35-44 | 35 | 65 | 50 |
| 45-54 | 30 | 70 | 50 |
| 55-64 | 25 | 75 | 50 |
| 65-74 | 20 | 80 | 50 |
| 75-84 | 15 | 85 | 50 |
| 85+ | 10 | 90 | 50 |

- Flow

- MaxPool
3x3+2(S)



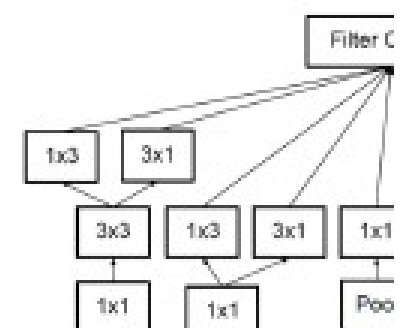
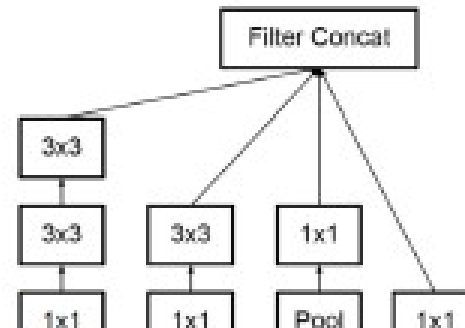
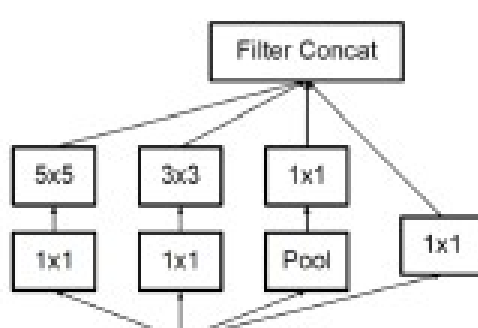
Conv NNs

□ Modern CNNs → GoogleLeNet / InceptionNet(2014)

■ Versions (v1, v2, v3, v4)

□ More templates, but the same 3 main properties are kept:

- Multiple branches
- Shortcuts (1x1, concat.)
- Bottleneck



Conv NNs

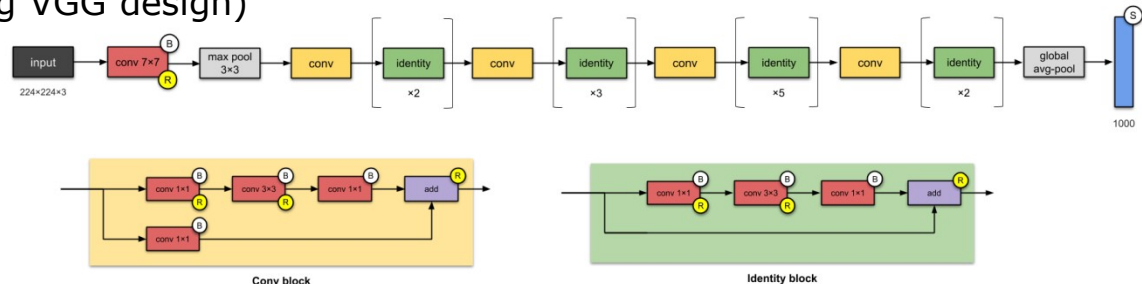
□ Modern CNNs → ResNet(2014) - Residual Neural Networks

■ Parents

- He et al. (Microsoft)
- He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015a. "Deep Residual Learning for Image Recognition." *CoRR* <http://arxiv.org/abs/1512.03385>
- Winner ILSVRC 2015

■ Flow

- 152 layers (following VGG design)



■ Particularities

- Batch normalisation after each convolution and before activation
- Ultra deep networks with residual connections
- Large learning rate (initial) 0.1
- No FC layers at the end

■ Input

- RGB image 224 x 224 x 3

■ Activation

- ReLU, than

■ Loss

- Softmax (cross-entropy)

■ # parameters

Conv NNs

□ Modern CNNs → ResNet(2014) - Residual Neural Networks

■ Flow

- 152 layers (following VGG design)

■ Particularities

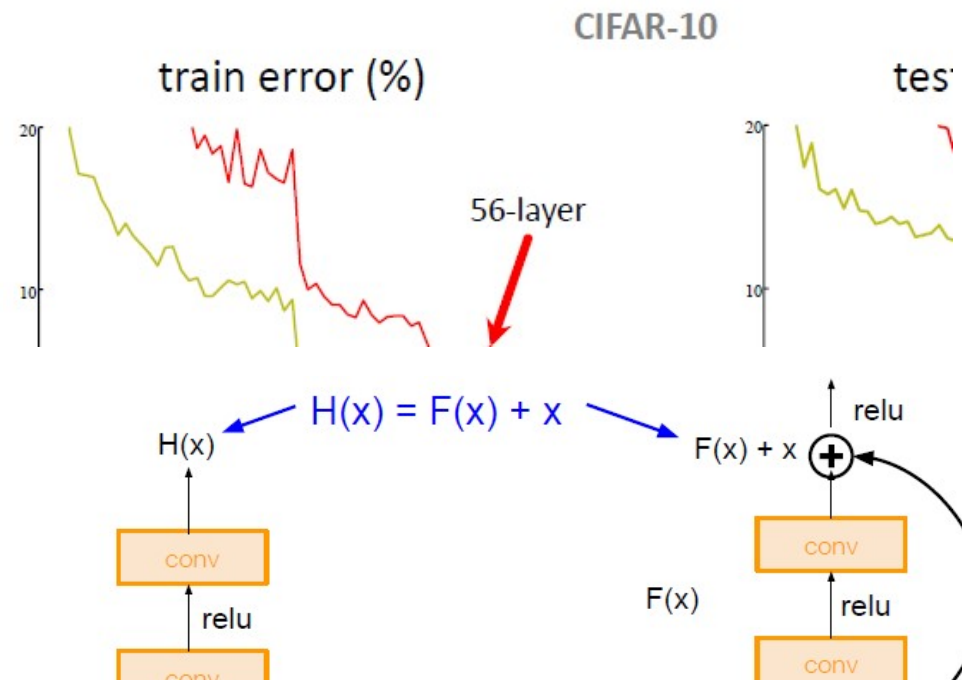
- Residual connections

■ Why?

- error information propagating back tends to get more and more diffused by the time it gets to the initial layers → the weights in the initial few layers are not modified in an optimal fashion → higher errors for deeper nets

■ How?

- Skip connections



Conv NNs

□ Modern CNNs → ResNet(2014) - Residual Neural Networks

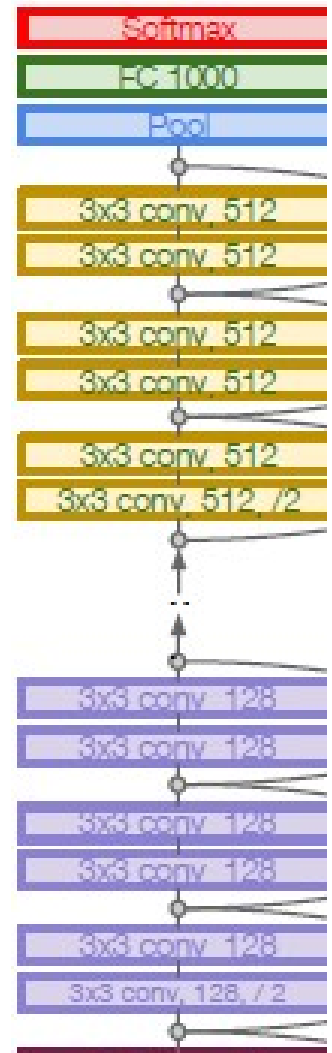
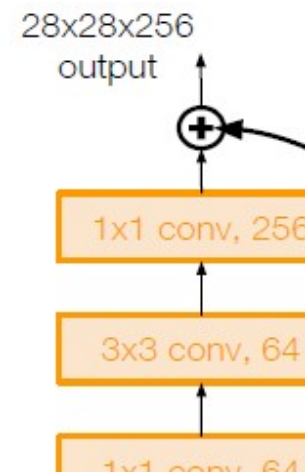
■ Flow

- 152 layers (following VGG design)

■ Particularities

- Residual connections
- A new layer = residual block (skip / shortcut connections)
 - 2 x (x 3 conv → batch norm → ReLU)
 - Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension)
 - ResNet ~ ensemble of shallow networks (paths of different lengths)
 - Help gradients to penetrate deeper into the network

- Additional conv layer at the beginning
- No FC layers at the end



Conv NNs

□ Modern CNNs → ResNet(2014) - Residual Neural Networks

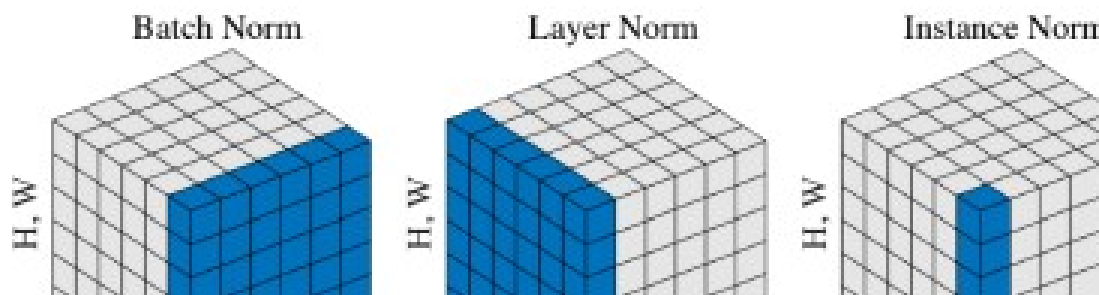
■ Flow

- 152 layers (following VGG design)

■ Particularities

- Batch normalisation after each convolution and before activation
 - Solves the problem of internal covariate shift (the distribution of each layer's inputs changes during training, as the parameters of the previous layers change)
 - Stabilise and accelerates the convergence
 - Group normalisation
 - No dropout layers

$$\text{BN}(\mathbf{x}) = \gamma \odot \frac{\mathbf{x} - \mu}{\sigma}$$



Input: Values of x over a mini-batch: \mathcal{B}

Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$$

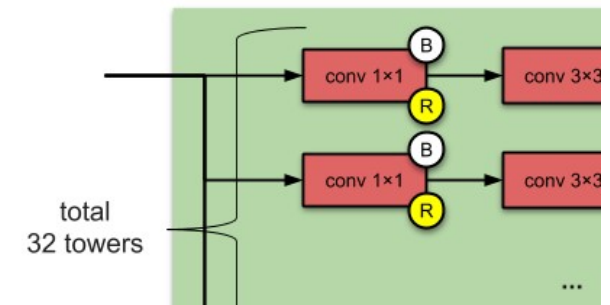
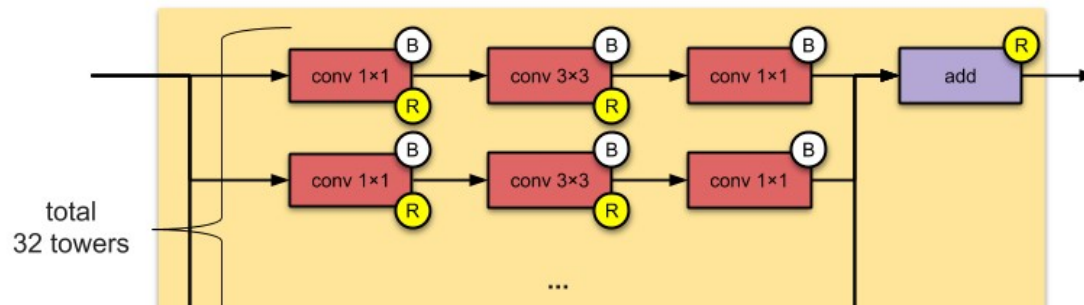
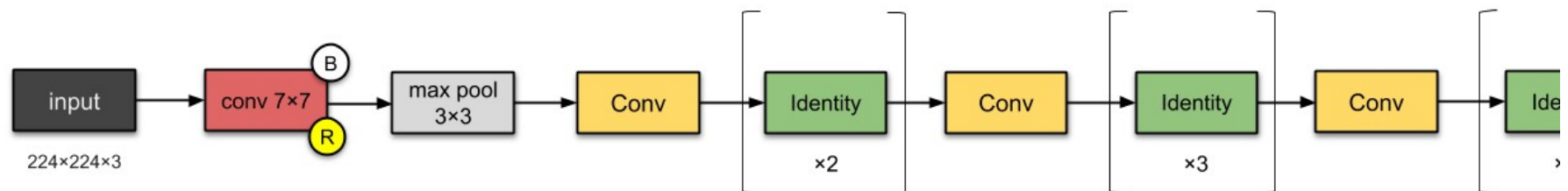
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ r}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sigma_{\mathcal{B}}}$$

- Large learning rate (initial) 0.1
 - Divided by 10 when the validation error plateaus
- No FC layers at the end

Conv NNs

- Modern CNNs → ResNeXt (2017)
 - Hybridisation of Inception and ResNet
 - <https://arxiv.org/pdf/1611.05431.pdf>

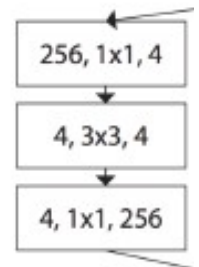
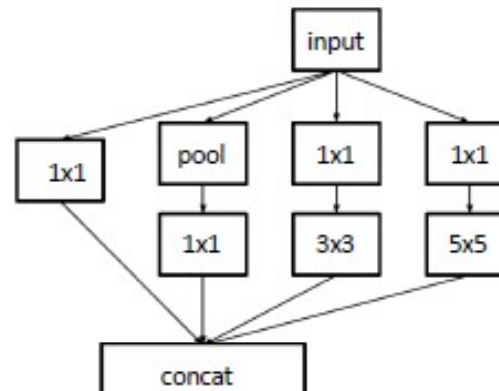


- Particularities
 - Shortcut
 - Bottleneck
 - Multi-branch

Conv NNs

Modern CNNs → ResNeXt (2017)

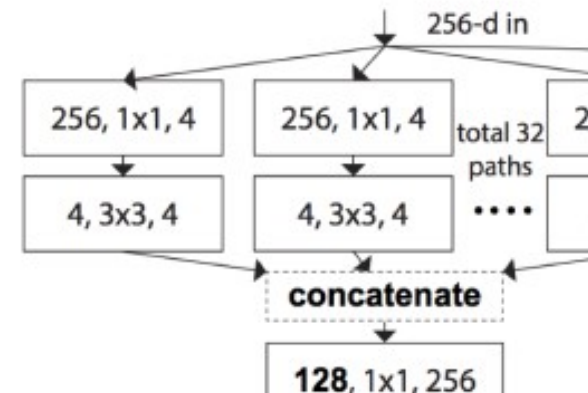
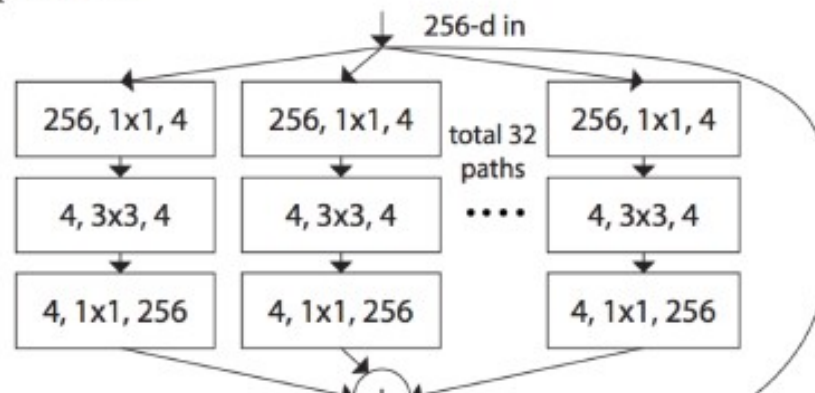
- Hybridisation of Inception and ResNet
- Particularities
 - Shortcut
 - Bottleneck



Multi-branch

- concatenation and addition are interchangeable → General property for Deep CNNs
- uniform multi-branching can be done by group-conv

equivalent



Conv NNs

□ Other architectures

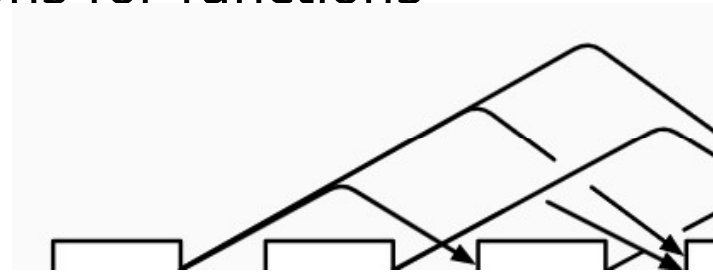
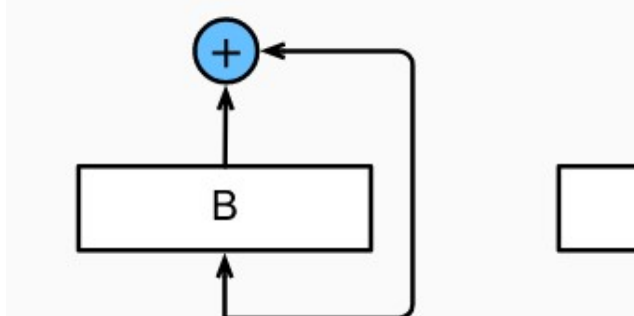
■ Inception-ResNet (2016)

■ Dense Net

□ Remember Taylor expansions for functions

$$f(x) = f(0) + f'(x)x + \frac{1}{2}f''(x)x^2 + \frac{1}{3!}f'''(x)x^3 + o(x^3)$$

$$f(x) = [x, f_1(x), f_2(x, f_1(x)), f_3(x, f_1(x), f_2(x, f_1(x))), \dots]$$



■ Xception, MobileNet (Google)

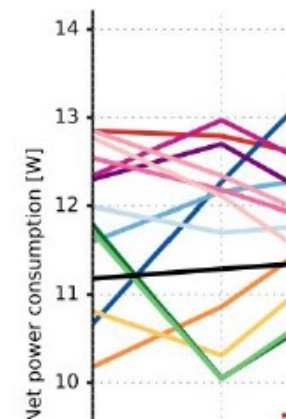
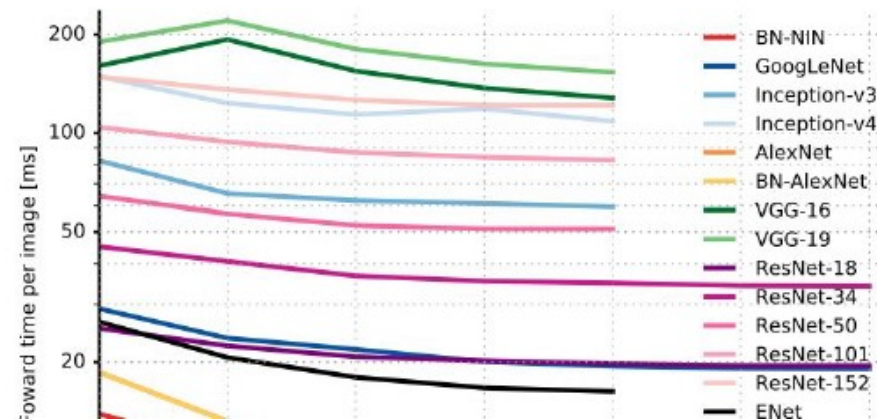
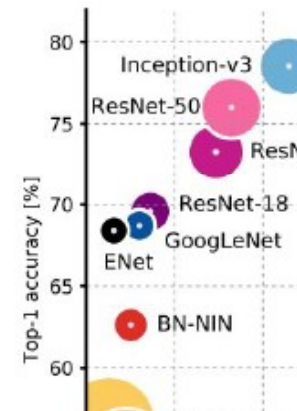
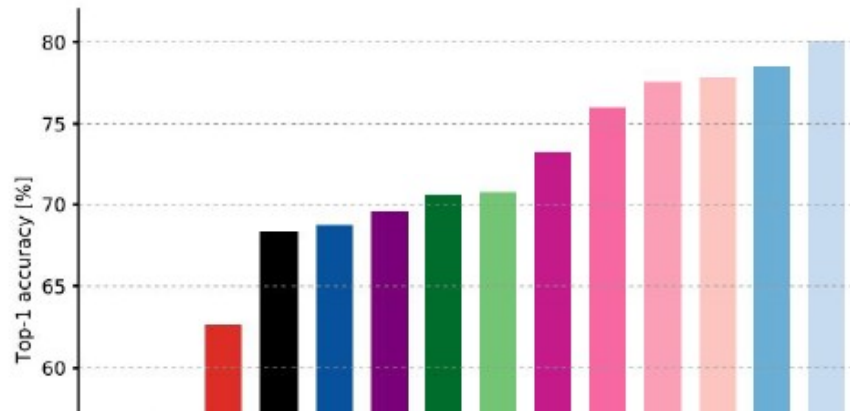
□ Depthwise convolutions (grouped conv with #channels= # group)

■ ...

Conv NNs

□ Review

- Complexity
- Forward pass time and power consumption



Conv NNs

□ Reducing overfitting

- increasing the amount of training data
 - Artificially expanding the training data
 - Rotations, adding noise,
- reduce the size of the network
 - Not recommended
- regularization techniques
 - Effect:
 - the network prefers to learn small weights, all other things being equal. Large weights will only be allowed if they considerably improve the first part of the cost function
 - a way of compromising between finding small weights and minimizing the original cost function (when λ is small we prefer to minimize the original cost function, but when λ is large we prefer small weights)
 - Give importance to all features
 - $X = [1, 1, 1, 1]$
 - $W_1 = [1, 0, 0, 0]$
 - $W_2 = [0.25, 0.25, 0.25, 0.25]$
 - $W_1^T X = W_2^T X = 1$
 - $L1(W_1) = 0.25 + 0.25 + 0.25 + 0.25 = 1$
 - $L1(W_2) = 1 + 0 + 0 + 0 = 1$

Conv NNs

□ Reducing overfitting - regularization techniques

■ Methods

- L1 regularisation – add the sum of the absolute values of the weights $C = C_0 + \lambda/n \sum |w|$
 - the weights shrink by a constant amount toward 0
 - Sparsity (feature selection – more weights are 0)
- *weight decay (L2 regularization)* - add an extra term to the cost function (the *L2 regularization term* = the sum of the squares of all the weights in the network = $\lambda/2n \sum w^2$): $C = C_0 + \lambda/2n \sum w^2$
 - the weights shrink by an amount which is proportional to w
- Elastic net regularisation
 - $\lambda_1|w| + \lambda_2 w^2$
- Max norm constraints (clapping)
- Dropout - modify the network itself (<http://www.cs.toronto.edu/~rsalakhu/papers/srivastava14a.pdf>)
 - Some neurons are temporarily deleted
 - propagate the input and backpropagate the result through the modified network
 - update the appropriate weights and biases.
 - repeat the process, first restoring the dropout neurons, then choosing a new random subset of hidden neurons to delete

Improve NN's performance

- ❑ Cost functions → loss functions
- ❑ Regularisation
- ❑ Initialisation of weights
- ❑ NN's hyper-parameters

Improve NN's performance

□ Cost functions → loss functions

■ Possible cost functions

□ Quadratic cost

- $\frac{1}{2n} \sum_x \|D - C\|^2$

□ Cross-entropy loss (negative log likelihood)

- $-\frac{1}{n} \sum_x [D \ln C + (1 - D) \ln(1 - C)]$

■ Optimizing the cost function

□ Stochastic gradient descent by backpropagation

□ Hessian technique

- Pro: it incorporates not just information about the gradient, but also information about how the gradient is changing

- Cons: the sheer size of the Hessian matrix

□ Momentum-based gradient descent

- Velocity & friction

Improve NN's performance

□ Initialisation of weights

■ Pitfall

- all zero initialization

■ Small random numbers

- $W = 0.01 * \text{random}(D, H)$

■ Calibrating the variances with $1/\sqrt{\#Inputs}$

- $w = \text{random}(\#Inputs) / \sqrt{\#Inputs}$

■ Sparse initialization

■ Initializing the biases

■ In practice

- $w = \text{random}(\#Inputs) * \sqrt{2.0/\#Inputs}$

Improve NN's performance

□ NN's hyper-parameters*

- Learning rate η
 - Constant rate
 - Not-constant rate
- Regularisation parameter λ
- Mini-batch size

*see Bengio's papers: <https://arxiv.org/pdf/1206.5533v2.pdf> and <http://www.jmlr.org/papers/volume13/bergstra12a/bergstra12a.pdf> or Snock's paper <http://papers.nips.cc/paper/4522-practical-bayesian-optimization-of-machine-learning-algorithms.pdf> 2019

Improve NN's performance

□ NN's hyper-parameters

■ Learning rate η

- Constant rate

- Not-constant rate

 - Annealing the learning rate

 - Second order methods

 - Per-parameter adaptive learning rate methods

Improve NN's performance

- NN's hyper-parameters - Learning rate η
 - Not-constant rate
 - Annealing the learning rate
 - **Step decay**
 - Reduce the learning rate by some factor every few epochs
 - $\eta = \eta * \text{factor}$
 - Eg. $\eta = \eta * 0.5$ every 5 epochs
 - Eg. $\eta = \eta * 0.1$ every 20 epochs
 - Exponential decay
 - $a = a_0 \exp(-kt)$,
where a_0 , k are hyperparameters and t is the iteration number (but you can also use units of epochs).
 - $1/t$ decay
 - $a = a_0 / (1 + kt)$
where a_0 , k are hyperparameters and t is the iteration number.

Improve NN's performance

- NN's hyper-parameters - Learning rate η

- Not-constant rate

- Second order methods

- Newton's method (Hessian)

- *quasi-Newton* methods

- L- BGFS (Limited memory Broyden–Fletcher–Goldfarb–Shanno)

- https://static.googleusercontent.com/media/research.google.com/ro//archive/large_deep_networks_nips2012.pdf

- <https://arxiv.org/pdf/1311.2115.pdf>

Improve NN's performance

□ NN's hyper-parameters - Learning rate η

■ Not-constant rate

□ Per-parameter adaptive learning rate methods

■ Adagrad

- <http://www.jmlr.org/papers/volume12/duchi11a/duchi11a.pdf>

■ RMSprop

- http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf

■ Adam

- <https://arxiv.org/pdf/1412.6980.pdf>

Tools

- ❑ Keras
 - NN API
 - <https://keras.io/>
 - + Theano (machine learning library; multi-dim arrays)
<http://www.deeplearning.net/software/theano/>
http://www.iro.umontreal.ca/~lisa/pointeurs/theano_scipy2010.pdf
 - + TensorFlow (numerical computation) <https://www.tensorflow.org/>
- ❑ Pylearn2 <http://deeplearning.net/software/pylearn2/>
 - ML library
 - + Theano
- ❑ Torch <http://torch.ch/>
 - scientific computing framework
 - Multi-dim array
 - NN
 - GPU
- ❑ Caffe
 - deep learning framework
 - Berkley

□ Presented information was collected from various sources:

- <https://cs230.stanford.edu/>
- <http://cs231n.stanford.edu/>
- <https://d2l.ai/>
- <https://berkeley-deep-learning.github.io/cs294-131-s17/>
- <https://machinethink.net/>
- <https://cedar.buffalo.edu/~srihari/CSE676/>
- <http://karpathy.github.io/>