	UNIVERZITET „UNION“ <b>RAČUNARSKI FAKULTET</b> Knez Mihailova 6/VI 11000 BEOGRAD	Broj:
		Datum:
		11.09.2018.

UNIVERZITET UNION  
RAČUNARSKI FAKULTET  
BEOGRAD  
Računarske nauke

# Sigurnost modela dubokog učenja

DIPLOMSKI RAD

Kandidat: Stanko Kovačević  
Broj indeksa: RN 29/13

Mentor: prof. Dr Dragan Urošević

Beograd, 18.12.2016.

# Sadržaj

1. Apstrakt.....	3
2. Oznake .....	4
3. Osnove dubokog učenja.....	5
3.1. Istorija i razvoj modela .....	5
3.2. Matematički model mašinskog učenja .....	6
3.2.1. Osnovni gradivni blokovi dubokih modela .....	8
3.2.2. Statistička estimacija.....	11
3.2.3. Matematička optimizacija.....	13
3.3. Duboki modeli – neuralne mreže .....	17
3.3.1. Konvolucione neuralne mreže.....	24
4. Napadi na modele dubokog učenja .....	24
4.1. Tehnike izvođenja napada.....	27
4.2. Tehnike odbrane .....	35
5. Zaključak.....	37
6. Reference .....	38

# 1. Apstrakt

Tema diplomskog rada je istraživanje i razvoj napadačkih i odbrambenih tehnika (eng. *Adversarial Attacks*) nad modelima dubokog učenja (eng. *Deep Learning Models*).

Veliki deo današnje tehnologije zasniva se na sistemima koji su sposobni da donesu odluke, i reše instance problema za koje prethodno nisu eksplicitno programirani. Tipičan primer je prepoznavanje objekta na slici. Ne postoji tačno eksplicitno pravilo prema kome možemo da napišemo egzaktno deterministički algoritam koji može da klasifikuje sadržaj slike. Čak i kada bi bilo moguće razviti takvo pravilo, postojao bi drugi praktični problem zbog toga što je problem klasifikacije šablona NP-hard. Danas se ove praktične poteškoće prevazilaze upotrebom mašinskog učenja (eng. *Machine Learning*), a najpopularniji modeli su duboke veštačke neuralne mreže (eng. *Deep Neural Networks*), koje pomoću statističke estimacije i matematičke optimizacije omogućavaju kreiranje algoritama koji iz velikog skupa podataka mogu sami da nauče *fuzzy* pravila za rešavanje problema koji su veoma značajni u praksi.

Problem sa pristupom mašinskog učenja je što skoro svi modeli uče distribuciju podataka, tj. uče karakteristike stohastičkog procesa koji generiše podatke, pa nad tim karakteristikama grade *fuzzy* semantičku reprezentaciju, a ne semantički model znanja direktno (npr. model će da nauči da ako se na slici pojavljuju nos, usta, oči, i uši zajedno, da je to slika lica sa verovatnoćom od npr. 0.99, ali model neće da nauči da lice ima usta, oči, uši, nos, itd). Zbog toga je moguće modifikovati ulaz u ove modele kako bi se distribucija iskrivila i tako omeo rad modela i model doterao u stanje koje odgovara napadaču. Ovo predstavlja veliki problem zato što su inteligentni sistemi svuda oko nas (banke, automobili bez vozača, medicina, sigurnost računarskih sistema, itd).

Cilj diplomskog rada jeste predstavljanje nekih od osnovnih ali veoma moćnih tehnika napada na modele dubokog učenja zasnovane na metodi gradijenta poput klasične propagacije unazad po slici, FGSM (eng. *Fast Gradient Sign Method*). Diplomski rad će također sadržati i tehnike za odbranu od ovih napada, i pokazati koliko je ovo težak problem. Tehnike za odbranu od ovakvih napada su otvoren problem u računarskim naukama, i ne postoji metod koji povoljno rešava i štiti od bilo kog napada. Također će biti prikazano kako se praktično izvodi ovaj napad pomoću aproksimacije modela, i svojstva prenosa protivničkih podprostora (eng. *Adversarial Subspace Transfer Property*).

Diplomski rad će sadržati i aplikacije napisane u programskom jeziku Pajton, koje će služiti kao primeri teorije koja se bude izlagala. U tu svrhu će da koriste biblioteke *TensorFlow* (biblioteka koja omogućava lako i jednostavno konstruisanje dubokih modela koji se pokreću na grafičkoj kartici), kao i raznih statističkih i grafičkih biblioteka za generisanje izveštaja i prikaza potrebnih za ilustracije ideja.

## 2. Oznake

Skup: Nakrivljeno veliko slovo, npr.  $A$

Skalar: Nakrivljeno malo slovo, npr.  $x, y, z$

Vektor: Nakrivljeno podebljano malo slovo, npr.  $\mathbf{x}, \mathbf{y}$ .  $i$ -ti element vektora  $\mathbf{x}$  se označava sa podskriptom, npr.  $x_i$

Matrica: Nakrivljeno podebljano veliko slovo, npr.  $\mathbf{\Theta}, \mathbf{X}, \mathbf{Y}$ . Elementi matrice se označavaju duplom podskriptom. Element u  $i$ -tom redu i  $j$ -toj koloni matrice  $\mathbf{X}$  se označava sa  $X_{i,j}$

Vektorski prostor matrica: Skup  $\mathcal{M}(n, m)$  predstavlja skup svih matrica koje imaju  $n$  redova i  $m$  kolona.

Tensor: Podebljano veliko slovo, npr.  $\mathbf{\Theta}$ . U ovom radu se koriste tensori 3 reda, tj. niz matrica.  $l$ -ta matrica tensora  $\mathbf{\Theta}$  se označava sa  $\mathbf{\Theta}^l$ , a konkretan element tensora u  $i$ -tom redu,  $j$ -toj koloni matrice  $l$  sa  $\Theta_{i,j}^l$

Skalarni (unutrašnji) proizvod: Ako su  $\mathbf{x}, \mathbf{y}$  vektori iz istog vektorskog prostora, njihov skalarni proizvod se definiše kao  $r = \langle \mathbf{x}, \mathbf{y} \rangle$ , gde  $r = \sum_i x_i * y_i$

Tensor (vanjski) proizvod: Ako su  $\mathbf{x}, \mathbf{y}$  vektori, njihov vanjski proizvod je matrica  $\mathbf{Z} = \mathbf{x} \otimes \mathbf{y}$  gde  $Z_{i,j} = x_i * y_j$

Hadamardov proizvod: Ako su  $\mathbf{x}, \mathbf{y}$  vektori, njihov Hadamardov proizvod se definise kao  $\mathbf{z} = \mathbf{x} \circ \mathbf{y}$  gde  $z_i = x_i * y_i$

Funkcija: Oznake funkcija se uvode na mestu na kome se prvi put pojavljuju. Funkcije mogu biti parametrizovane, a parametri se postavljaju posle znaka  $;$ . Parametri su ulaz u funkciju koji je fiksiran, npr.  $f(\mathbf{x}; \mathbf{\Theta})$ .

Uslovna verovatnoća: U tekstu ćemo koristiti  $p(\mathbf{y}|\mathbf{x})$  da označimo uslovnu verovatnoću realizacije vektora  $\mathbf{y}$  ako znamo da se realizovao vektor  $\mathbf{x}$

### 3. Osnove dubokog učenja

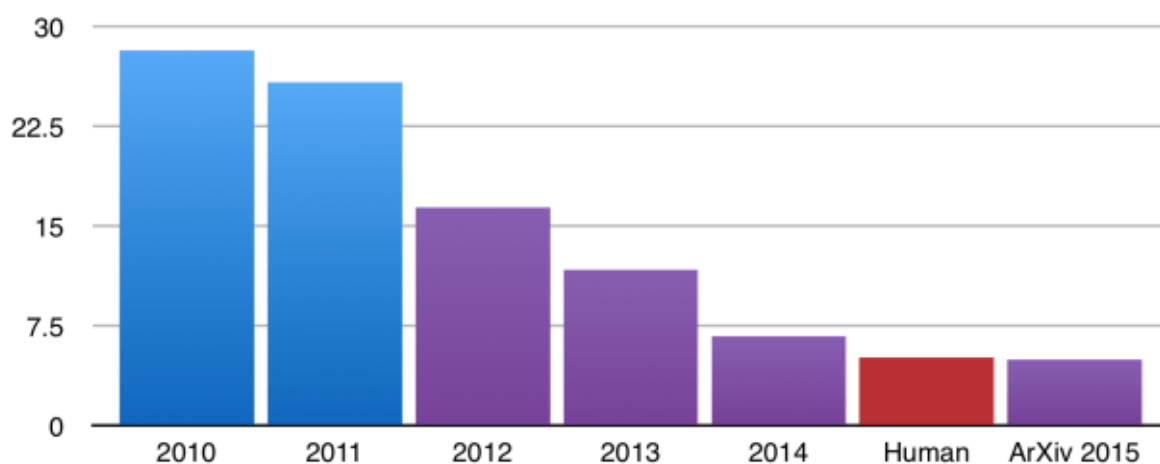
Još od pojave prvih mašina, ljudi su bili opsjednuti time da nateraju mašine da budu autonomne, i da razmišljaju kao ljudi. Razlog tome je ljudska priroda da bude radoznao, ali i velika praktična primena u industriji i poboljšavanju kvaliteta života ljudi. Međutim, ispostavilo se da je ovakve mašine veoma teško napraviti iz dva prosta razloga. Prvi razlog je mehanički. Čovek ima veoma preciznu kontrolu nad svim svojim zglobovima, i veoma složen i precizan kinetički lanac, koji sa centralnim nervnim sistemom održava ravnotežu čoveka, koordiniše svim njegovim pokretima, pruža povratne informacije u vidu čula. Konstrukcija motora koji je precizan i energetski efikasan kao ljudski zglobovi je praktično zahtevna. Drugi, daleko složeniji problem je inteligencija. Ljudski centralni nervni sistem je povezan sa svim sensorima (čulima), aktivatorima (kinetičkim lancem i zglobovima), i omogućava njihovu koordinaciju. Međutim, ljudski mozak omogućava i daleko složenije stvari kao što su donošenje odluka, učenje, apstrakcija okruženja, rekonstrukcija okruženja iz nervnih impulsa poslatih iz organa za vid i organa za sluh.

Uzmimo na primer ljudski vid. Čovek genetski ima ugrađene određene sposobnosti i karakteristike čula vida koje su se razvile procesom evolucije, kao na primer strah od zmija i paukova [1], fokusiranje na lica, kao i čitav sistem koji rekonstruiše mentalnu sliku onog što ljudske oči vide. Međutim, da bi čovek funkcionisao, i preživeo u svom okruženju, on mora konstantno da uči i da povezuje te mentalne slike sa prethodnim iskustvom. Na primer, da bi čovek mogao da čita knjigu, on prvo mora da nauči kako knjiga izgleda. Ovaj proces izvlačenja informacije sa mentalne slike i generalizacija oblika knjige je veoma efikasan kod ljudi, i čovek veoma brzo generalizuje i uči da prepozna svakodnevne predmete.

#### 3.1. Istorija i razvoj modela

Može se reći da mašinsko učenje potiče još od 1763 godina pojavom matematičke statistike i početka Bajesove teoreme [2] na kojoj se zasniva većina algoritama mašinskog učenja, ali tek od 1930-tih godina se ova disciplina aktivno izučava kao oblast veštačke inteligencije. Istraživanja u neurologiji, Šenonova teorija informacija, i Alan Turingova teorija izračunljivosti su doprinele razvoju ideja prvih agenata veštačke inteligencije. MekKulok, i Pits [3] su prvi koji su pokušali da naprave model koji će kasnije postati popularan kao neuralne mreže. 1950 godine Alan Turing postavlja osnove veštačke inteligencije Turingovim testom [4], a već 1951 godine Marvin Minski konstruiše prvi neuralni računar *SNARC*. Svi dosadašnji modeli nisu imali sposobnost učenja, nego su samo opisivali algoritamski potrebne korake pod uslovom da su svi parametri već zadani. To se menja 1957 godine kada Frenk Rosenblatt objavljuje algoritam učenja Perceptron. Iako su na početku kreirali veliko interesovanje u akademiji, 1970 godina mašinsko učenje ulazi u svoje mračno doba. Minski je pokazao da neuralne mreže imaju ograničenja, i da Perceptron ne može da nauči probleme koji su linearni, ali da ne može da nauči složenije funkcije poput XOR. Jedan od najznačajnijih događaja koji neuralne mreže vraća u život jeste razvoj algoritma učenja sa propagacijom

unazad [5], koji omogućava da se treniraju složeni neuralni modeli sposobni da nauče da aproksimiraju nelinearne funkcije. Ovaj algoritam se u različitim oblicima koristi i danas. U istom periodu se pojavljuju i razne arhitekture neuralnih mreža kao što su rekurentne neuralne mreže (Hopfield) [6], kao i konvolucione neuralne mreže (*Li Kun*) [7], a kasnije 1997 se pojavljuje i LSTM arhitektura rekurentnih neuralnih mreža (Hošreiter, Smithaber) [8]. Iako



Ilustracija 1: top-5 greška na ImageNet skupu

su sve tehnike koje se koriste i danas razvijene mnogo ranije, ipak sve do 2012 godine neuralne mreže su bile potisnute od strane statističkih algoritama poput mašina vektora podrške (eng. *Support Vector Machines, SVM*), stabala odlučivanja, i slično. Razlog maloj popularnosti neuralnih mreža je bila složenost treninga, jer je za uspešnu konvergenciju potreban veliki skup trening podataka, kao i velika procesorska moć. Ranije nije bilo moguće prikupiti dovoljnu količinu podataka, ali poboljšanjem pristupačnosti tehnologije, postaje lakše prikupiti podatke. Također, tih godina počinje era primene grafičkih kartica u naučne svrhe. 2012 godina je prekretnica dubokog učenja kada na *ImageNet* takmičenju po prvi put pobeđuje duboka neuralna mreža *AlexNet* [9].

### 3.2. Matematički model mašinskog učenja

Mašinsko učenje se može definisati na razne načine, međutim, opšte prihvaćena definicija je sledeća:

---

**Definicija 1:** Za program kažemo da uči da reši zadatak  $T$  ako se sa povećanjem iskustva  $E$ , povećavaju i performanse  $P$  programa.

---

Primer zadatka bi bio klasifikacija da li je na slici mačka ili ne, ili predikcija cene kuće na osnovu parametara kao što su površina, lokacija, starost kuće. Ovakvo učenje se naziva učenje nadgledanjem (eng. *Supervised Learning*). Iskustvo uglavnom predstavlja skup podataka koji se izvlači i prikuplja iz populacije, ili se simulira i meri kod učenja pojačavanjem (eng. *Reinforcement learning*). Primer skupa podataka za trening klasifikatora mačke bi bio

skup slika od polovina sadrži mačku, a druga polovina ne. Postoji čitava oblast o tome kako se prikupljaju podaci za test i trening, i u nastavku se nećemo baviti tim problemima i pretpostavićemo da su podaci iz trening skupa statistički nezavisni, i da imaju istu statističku raspodelu verovatnoće. Performanse programa predstavljaju meru koliko model dobro aproksimira podatke koje je video u toku treninga, tj. Koliko dobro aproksimira funkciju koja mu je zadana da je nauči. Konkretni detalji funkcije performanse zavise od problema koji se rešava, a najčešće je to srednja kvadratna greška (eng. *Mean Squared Error*, MSE, L2), preciznost, odziv, ROC kriva (eng. *Receiver Operator Characteristic*), i druge. U ovom radu ćemo se fokusirati na problem klasifikacije tako da će nam performansa  $P$  biti preciznost i odziv koje ćemo definisati kasnije. Sve tehnike opisane u nastavku ovog rada analogno se mogu primeniti i na druge funkcije  $P$  pod uslovom da su te funkcije diferencijalne skoro svuda. Ovakvo ograničavanje na specifičan problem klasifikacije omogućava da se iznesu ciljne ideje, i da se prije toga čitalac uvede u materiju bez toga da se mu se odvuče pažnja od primarnog cilja ovog rada, a to su napadi na duboke modele. Zadatak  $T$  će u daljem tekstu biti zadatak klasifikacije, tj. mapiranje ulaznog signala u vektor raspodele verovatnoća

$$\mathbf{y}' = f(\mathbf{x}; \boldsymbol{\theta}), y'_j > 0, \sum_j y'_j = 1$$

tako da se aproksimira stohastički proces  $\mathbf{y} = \tau(\mathbf{x})$  koji generiše podatke, i koji nam je nepoznat. Koncept vektor  $\mathbf{y}$  ima kategoričnu raspodelu, gde  $i$ -ti element ima vrednost 1 ako se u ulaznom signalu pojavljuje koncept  $i$ , i on je observabilan.  $f$  je familija funkcija kojom pokušavamo da aproksimiramo nepoznati model  $\tau$ , a  $\boldsymbol{\theta}$  je parametar koji definiše konkretnu funkciju iz familije  $f$ . Preciznije, koncept vektor definišemo na sledeći način.

---

**Definicija 2:** Neka je  $\mathbf{y}$  diskretan  $m$  dimenzionalni slučajni vektor, čija podrška je

$$C_m = \left\{ \mathbf{y} \in \{0, 1\}^m : \sum_{j=1}^m y_j = 1 \right\}$$

Neka su  $p_1, p_2, \dots, p_m$  pozitivni brojevi tako da

$$\sum_{j=1}^m p_j = 1$$

Kažemo da je  $\mathbf{y}$  koncept vektor ako je udružena marginalna raspodela:

$$p(y_1, y_2, \dots, y_m) = \begin{cases} \prod_{j=1}^m p_j^{y_j}, & (y_1, y_2, \dots, y_m) \in C_m \\ 0, & (y_1, y_2, \dots, y_m) \notin C_m \end{cases}$$

---

Primer: Neka imamo model  $\mathbf{y}' = f(\mathbf{x}; \boldsymbol{\theta})$  koji klasifikuje da li se na slici nalazi mačka, pas, ili ptica. Vektor  $\mathbf{x}$  bi u ovom slučaju predstavljao sliku, dok bi se koncept vektor mogao definisati na sledeći način:  $\mathbf{y} = (1, 0, 0)$  ako je na slici mačka,  $\mathbf{y} = (0, 1, 0)$  ako je na slici pas,

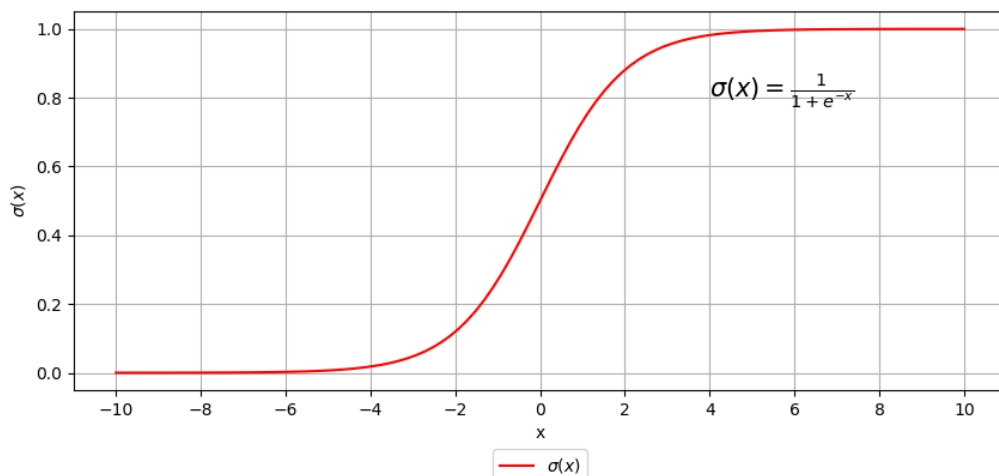
ili  $\mathbf{y} = (0, 0, 1)$  ako je na slici ptica. Izlaz  $\mathbf{y}'$  bi bio vektor koji aproksimira uslovnu raspodelu  $p_f(\mathbf{y}|\mathbf{x})$ . Iako je marginalna raspodela korisna kod eksploratorne analize, nas će više zanimati uslovna raspodela  $p_h(\mathbf{y}|\mathbf{x})$ , gde je  $h$  model koji opisuje generativni proces. U slučaju kada je  $h = \tau$ , tj. kod potpunog opisa,  $p_h(\mathbf{y}|\mathbf{x}) = \mathbf{y}$ , tj. za svaki ulazni vektor se tačno zna koji je koncept (pretpostavka je da je proces deterministički).

### 3.2.1. Osnovni gradivni blokovi dubokih modela

Glavno pitanje kod dizajna sistema zasnovanog na mašinskom učenju jeste koju funkciju izabrati za aproksimaciju. Za početak možemo se ograničiti na konkretnu familiju funkcija  $f$ , kao što je logistička funkcija

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Ova funkcija modeluje verovatnoću  $p_f(y = 1|\mathbf{x})$ , gde je  $y$  skalarna vrednost, a pošto direktno modeluje raspodelu, spada u diskriminacione modele. U slučaju binarne klasifikacije koncept vektor  $\mathbf{y}$  se obično modeluje kao skalarna vrednost i to  $y = \begin{cases} 1, & \mathbf{y} = (1, 0) \\ 0, & \mathbf{y} = (0, 1) \end{cases}$



Ilustracija 2: Sigmoidalna funkcija

Kada uvrstimo ovu funkciju, naš model postaje

$$\mathbf{y}' = f(\mathbf{x}; \boldsymbol{\theta}) = \sigma(b + \langle \mathbf{x}, \boldsymbol{\theta} \rangle) = \frac{1}{1 + e^{-(b + x_1 * \theta_1 + x_2 * \theta_2 + \dots + x_k * \theta_k)}}, \quad f: \mathbb{R}^k \rightarrow (0, 1)$$

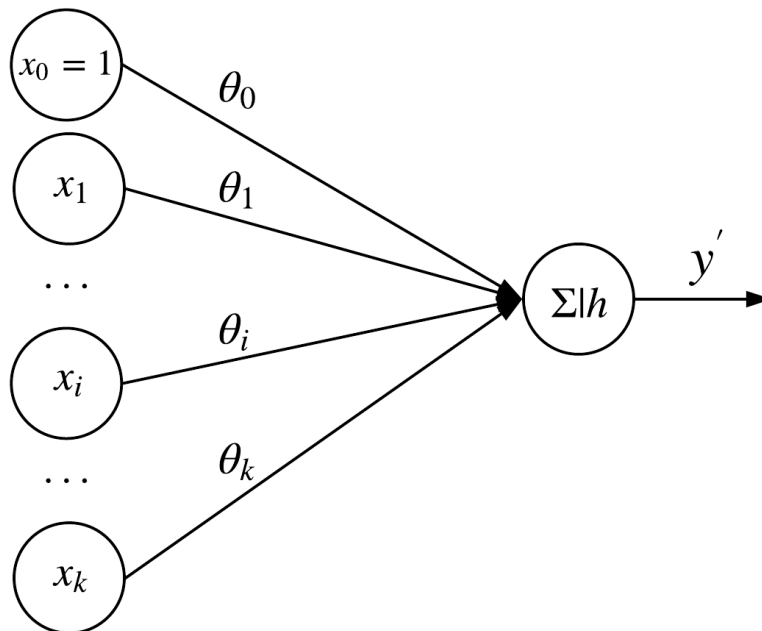
Ovim modelom modelujemo uslovnu verovatnoću da je na slici koncept 1, ako je ulaz  $\mathbf{x}$ .

Ako stavimo da je  $x_0 = 1$ ,  $b = \theta_0$ , tada dobijamo kompaktniji način reprezentacije koji se može vektorizovati i efikasno izračunati na grafičkom procesoru. U nastavku ćemo



pretpostavljati da je odrađeno pre procesiranje i da se ulaznom vektoru  $x$ , pridružio element  $x_0 = 1$ , i da je  $\theta$   $k+1$  dimenzionalni vektor parametara, pa logistička funkcija poprima oblik

$$f(x; \theta) = \sigma(\langle x, \theta \rangle) = \frac{1}{1 + e^{-\langle x, \theta \rangle}} = \frac{1}{1 + e^{-(x_0 * \theta_0 + x_1 * \theta_1 + x_2 * \theta_2 + \dots + x_k * \theta_k)}}, f: \mathbb{R}^{k+1} \rightarrow (0, 1)$$

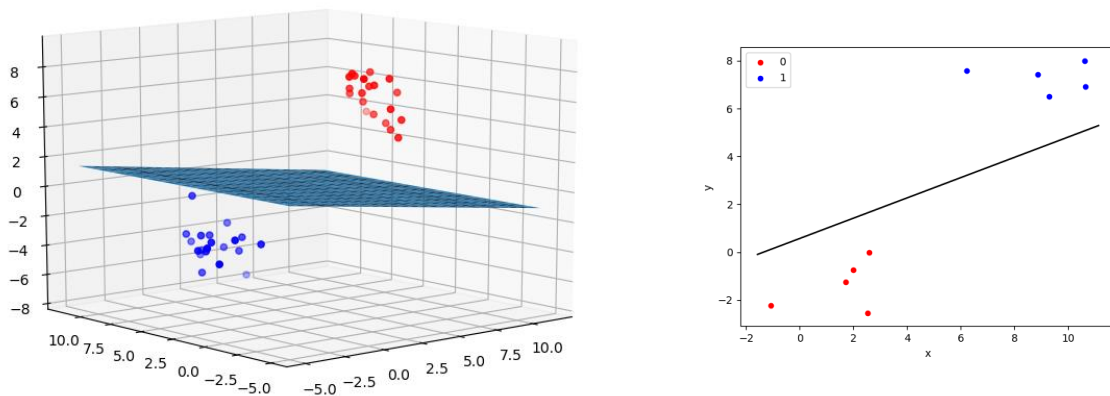


Ilustracija 3: Graf računa logističke regresije

Elementi ulaznog vektora  $x$  se množe sa elementima vektora parametara  $\theta$ , zatim se sabiraju, i propuštaju kroz aktivacionu funkciju  $h$ , koja je u ovom slučaju  $\sigma$ . Ovaj blok iznad je osnovni gradivni blok neuralnih mreža kao što ćemo videti kasnije.

Ova klasa funkcija deli  $k + 2$  dimenzionalni prostor na 2 dela sa  $k + 1$  dimenzionalnom hiper-ravni ( $k$  dimenzionalan podprostor)

$$\langle x, \theta \rangle = \sum_{j=0}^k \theta_j x_j = 0$$



Ilustracija 4: (Levo) Granica podele u 2d ravni. Sve tačke iznad linije sa klasifikuju kao plave, a ispod linije kao crvene. (Desno) Granica podele u 3d prostoru.

Logistička regresija je veoma bitan model zbog toga što predstavlja diferencijalnu verziju Hevisajdove funkcije koja je ključna kod modela neuralnih mreža, i predstavlja gradivni blok klasičnih neuralnih mreža. Logistička regresija se generalizuje na više kategoričnu klasifikaciju pomoću *softmax* funkcije

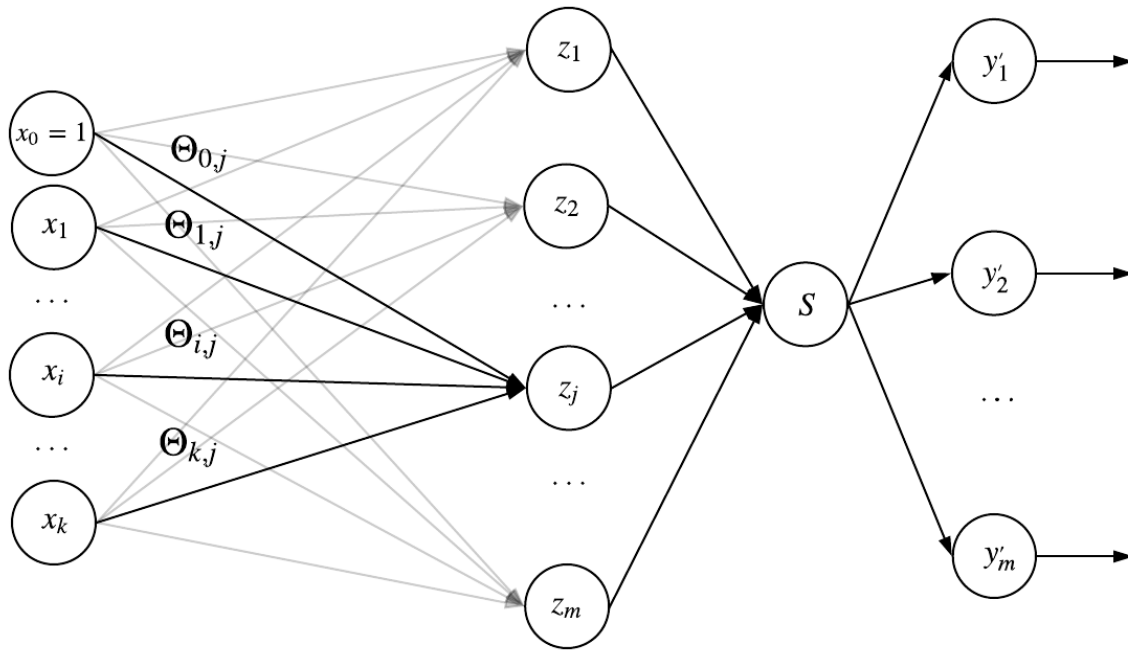
$$\mathbf{s} = S(\mathbf{z}), s_j = \frac{e^{z_j}}{\sum_{p=1}^k e^{z_p}}, s_j > 0, \sum_{j=1}^m s_j = 1, S: \mathbb{R}^k \rightarrow (0, 1)^m$$

*Softmax* funkcija sabija vektor  $\mathbf{z}$  u vektor verovatnoća i tako omogućava konstrukciju diskriminacionog modela  $p_f(\mathbf{y}|\mathbf{x})$ , gde je  $\mathbf{x}$  ulazni vektor, a  $\mathbf{y}$  koncept vektor. Umesto vektora parametara  $\boldsymbol{\theta}$ , sada imamo matricu parametara  $\boldsymbol{\Theta}$ , gde svaka kolona  $\boldsymbol{\Theta}_{:,i}$  predstavlja ugrađeni vektor parametara za i-ti koncept.

$$\mathbf{y}' = f(\mathbf{x}; \boldsymbol{\Theta}) = S(\mathbf{x}^T \boldsymbol{\Theta}),$$

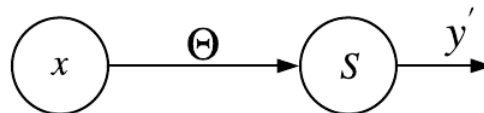
$$f: \mathbb{R}^{k+1} \rightarrow (0, 1)^m$$

$$\boldsymbol{\Theta} \in \mathcal{M}(k+1, m)$$



Ilustracija 5: Graf softmax regresije.

Model *softmax* regresije je sposoban da klasifikuje ulazni signal kad postoje više od 2 koncepta. Drugi način da ilustrujemo *softmax* regresiju jeste da grupišemo elemente vektora u jedan čvor. Ovaj način će biti veoma koristan kasnije kada pređemo na duboke modele jer bi u suprotnom bilo veoma teško predstaviti toliko veliki broj čvorova.



Ilustracija 6: Kompaktan prikaz softmax regresije

### 3.2.2. Statistička estimacija

Postavlja se pitanje kako izabrati parametre  $\theta$  da bi model konvergirao i rezultirao dobrim performansama, tj. da bi model aproksimirao funkciju koja mu je data da nauči, a u našem slučaju, to je raspodela  $p_{\tau}(\mathbf{y}|\mathbf{x})$ . Za podatke koje smo već videli u trening skupu podataka tačno znamo raspodelu  $p_{\tau}(\mathbf{y}|\mathbf{x})$ , međutim, naš cilj je da naučimo raspodelu za čitavu populaciju, a ne samo mali uzorak koji nam je potpuno poznat. Kao što smo već rekli, uslovnu raspodelu  $\mathbf{y} = p_{\tau}(\mathbf{y}|\mathbf{x})$  nad populacijom ćemo aproksimirati raspodelom

$$\mathbf{y}' = p_f(\mathbf{y}|\mathbf{x}) = f(\mathbf{x}; \theta)$$

Da bi našli parametre  $\Theta$ , potrebno je da u svakom trenutku znamo koliko je naša trenutna aproksimacija dobra, tj. koliko je  $p_\tau$  daleko od  $p_f$ . Drugim rečima, potrebna nam je metrika, tj. distanca između 2 slučajne raspodele. Jedna takva metrika (polu-metrika) je unakrsna entropija [10]

---

**Definicija 3:** Unakrsna entropija između dve slučajne promenjive  $p$ , i  $q$  se definiše kao  $H(p, q) = E_p[-\log q] = -\sum_x p(x) \log q(x)$

---

Unakrsna entropija ima nekoliko interpretacija. Prva zasnovana na teoriji informacija pokazuje koliko bita u proseku treba da se kodiraju simboli izvora ako se koristi šema kodiranja optimalna za stohastički izvor sa distribucijom  $q$  koja može biti procenjena nekom metodom, umesto prave distribucije  $p$  koja je nepoznata. Tačnije, ovo je entropija koja pokazuje prosečan broj bitova potreban ako se koristi „pogrešna“ distribucija.

---

**Teorema 1:** Unakrsna entropija  $H(p, q)$  je uvek veća od entropije  $H(p)$ , ili formalno:  $H(p, q) > H(p) \quad \forall p, q \in \Omega$ , gde je  $\Omega$  skup svih raspodela verovatnoće.

---

Dokaz: Šenon [11] je pokazao da prosečan broj bita potreban da se bez gubitka kodira izvor ne može biti manji od njegove entropije  $H(p)$ , gde je  $p$  parametar stohastičkog izvora koji predstavlja raspodelu verovatnoće pojavljivanja simbola. Dakle, ako je naša procenjena distribucija „pogrešna“, kodiranje neće biti optimalno, tj.  $H(p, q) > H(p)$ .

Unakrsna entropija nam omogućava da izmerimo rastojanje između 2 slučajna vektora  $\mathbf{y}$ , i  $\mathbf{y}'$ , tj. za svaki primerak iz trening skupa  $\mathbf{x}$ , možemo da izmerimo koliko je naš model daleko od stvarnog. Ta mera se naziva *loss* funkcija

$$L(\mathbf{y}, \mathbf{y}') = L(p_\tau(\mathbf{y}|\mathbf{x}), p_f(\mathbf{y}|\mathbf{x})) = H(p_\tau(\mathbf{y}|\mathbf{x}), p_f(\mathbf{y}|\mathbf{x}))$$

$$L(\mathbf{y}, \mathbf{y}') = H(\mathbf{y}, \mathbf{y}') = -\sum_{i=1}^m y_i \log y'_i$$

Poslednji deo slagalice jeste da se *loss* funkcija minimizuje na čitavom trening skupu u proseku, a ne samo za jedan ulaz.

---

**Definicija 4:** Funkcija greške modela  $\mathbf{y}' = f(\mathbf{x}; \Theta)$ , na trening skupu  $E$  se definiše na sledeći način:

$$C(\Theta) = \frac{1}{|E|} \sum_{(\mathbf{x}, \mathbf{y}) \in E} L(\mathbf{y}, f(\mathbf{x}; \Theta))$$


---

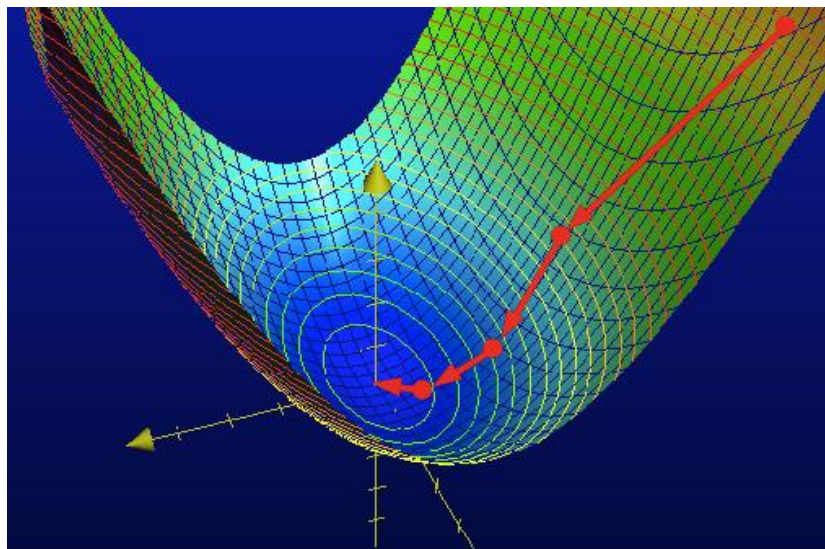
Pošto trening skup  $E$  predstavlja slučajan uzorak iz neke populacije, svaka estimacija treba da se generalizuje na čitavu populaciju, tj. pristrasnost  $bias(f) = E(p_f) - p_\tau$  treba da bude što manja, kao varijansa modela  $Var(f) = Var(p_f)$ . Greška modela se može dekomponovati na

$$Err(f) = bias(f)^2 + Var(f) + \epsilon$$

gde je  $\epsilon$  ne merljiva greška. Nažalost, tačna vrednost svih izraza gore nije rešiva, i mogu se raditi samo neke aproksimacije. Generalno pravilo je da kompleksni model spušta pristrasnost ali podiže varijansu modela, a manje kompleksan model podiže pristrasnost, ali spušta varijansu. Glavni problem u praksi kod rešenja sa mašinskim učenjem jest pronalaženje neke „zlatne sredine“ u kojoj imamo dovoljno kompleksan model koji ima kapacitet da nauči raspodelu, ali ne previše tako da se loše generalizuje na populaciju [12, Ch. 7].

### 3.2.3. Matematička optimizacija

Kada smo tačno definisali funkciju greške nad prostorom parametara, sve što je ostalo jeste da se ta funkcija minimizuje. Analitičko rešenje je veoma teško izvesti, pa se primenjuju iterativni metodi zasnovani na gradijentu funkcije. S toga se zahteva da svi modeli, budu diferencijabilni. Optimizacija se vrši tako što se određen broj puta pomeramo u prostoru parametara u smeru suprotnom gradijentu [12, Ch. 4]. Razlog je taj što gradijent



Ilustracija 7: Spust gradijentom. Algoritam počinje u gornjem uglu i iterativno se spušta do lokalnog maksimuma

funkcije uvek pokazuje u smeru u kome je porast funkcije najveći, pa je očigledno da ako se u toj tački pomerimo suprotno od gradijenta, da imamo najmanji porast, što je u suštini najveći spust.

Spust gradijentom uvek pronalazi rešenje ukoliko je funkcija konveksna. Na žalost, kod dubokih neuralnih mreža, funkcija greške nije konveksna ali to u praksi ne predstavlja problem. Pri svakoj iteraciji, primenjuje se sledeća jednačina

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \nabla C$$

Ovo pravilo se ponavlja određen broj puta ili dok model ne konvergira.

Na primeru ćemo pokazati kako se trenira *softmax* klasifikator. Za to nam je potrebno da znamo  $\nabla C = (\frac{\partial C}{\partial \theta_{0,1}}, \frac{\partial C}{\partial \theta_{1,1}}, \dots, \frac{\partial C}{\partial \theta_{k,m}})$ .

$$C(\boldsymbol{\theta}) = \frac{1}{|E|} \sum_{(x,y) \in E} L(y, f(x; \boldsymbol{\theta}))$$

$$1. \frac{\partial C}{\partial \theta_{i,j}} = \frac{1}{|E|} \sum_{(x,y) \in E} \frac{\partial L}{\partial \theta_{i,j}}$$

$$2. \frac{\partial L}{\partial \theta_{i,j}} = \frac{\partial L}{\partial \mathbf{y}'} \frac{\partial \mathbf{y}'}{\partial \theta_{i,j}} = \frac{\partial L}{\partial \mathbf{y}'} \frac{\partial \mathbf{y}'}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \theta_{i,j}}$$

$$3. \frac{\partial L}{\partial \mathbf{y}'} = \nabla L, \text{ gde } \nabla L_k = \frac{\partial H(\mathbf{y}, \mathbf{y}')}{\partial y'_k} = -\frac{y_k}{y'_k}$$

$$4. \frac{\partial \mathbf{y}'}{\partial \mathbf{z}} = \frac{\partial S(\mathbf{z})}{\partial \mathbf{z}} = \mathbf{J}_S$$

$$\text{gde } \mathbf{z} = \mathbf{x}\boldsymbol{\theta} = (\mathbf{x}\boldsymbol{\theta}_{:,1}, \dots, \mathbf{x}\boldsymbol{\theta}_{:,m})$$

$$J_{S_{i,j}} = S(\mathbf{z})_i (\delta_{ij} - S(\mathbf{z})_j) = y'_i (\delta_{ij} - y'_j) - m \times m \text{ Džakobijeva matrica softmax funkcije}$$

$$\delta_{ij} = \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases}$$

$$5. \frac{\partial \mathbf{z}}{\partial \theta_{i,j}} = \frac{\partial (\mathbf{x}\boldsymbol{\theta})}{\partial \theta_{i,j}} = \mathbf{Z} = \mathbf{x} \otimes \mathbf{I}$$

gde je  $\mathbf{I}$  matrica identičnosti,  $\mathbf{Z}$  tensor 3 reda i definisan je kao  $Z_{l,j}^i = \delta_{l,j} x_i$

4,3 ->2:

$$6. \frac{\partial L}{\partial \theta_{i,j}} = (\nabla L J_s) \frac{\partial \mathbf{z}}{\partial \theta_{i,j}} = (\mathbf{y}' - \mathbf{y}) \frac{\partial \mathbf{z}}{\partial \theta_{i,j}}$$

5->6:

$$7. \frac{\partial L}{\partial \theta_{i,j}} = (\mathbf{x} \otimes (\mathbf{y}' - \mathbf{y}))_{i,j} = (y'_j - y_j)x_i$$

Na kraju kad uvrstimo 7 u 1, dobijamo konačnu jednačinu:

$$8. \frac{\partial C}{\partial \theta_{i,j}} = \frac{1}{|E|} \sum_{(x,y) \in E} (y'_j - y_j)x_i$$

U praksi se ipak ovakav oblik jednačine 8 ne koristi jer u ovakvom obliku gubimo prednosti izvršavanja algoritma na grafičkoj kartici. Kako je grafička kartica veoma dobra u množenju matrica, potrebno je sve algoritme vektorizovati.

Ako sa  $(\mathbf{x}^k, \mathbf{y}^k)$  označimo k-ti element iz trening skupa  $E$ , tada jednačina 8 postaje:

$$9. \frac{\partial C}{\partial \theta_{i,j}} = \frac{1}{|E|} \sum_{k=1}^{|E|} (y_j'^k - y_j^k)x_i^k$$

Ako sve elemente trening skupa, spakujemo u  $|E| \times k$  matricu  $\mathbf{X}$ , a odgovarajuće koncept vektore u  $|E| \times m$  matricu  $\mathbf{Y}$ , možemo jednačinu 9 zapisati kao:

$$10. \frac{\partial C}{\partial \theta_{i,j}} = \frac{1}{|E|} \sum_{k=1}^{|E|} (Y'_{k,j} - Y_{k,j})X_{k,i} = \frac{1}{|E|} \sum_{k=1}^{|E|} (\mathbf{Y}' - \mathbf{Y})_{k,i} X_{k,j}$$

Jednačina iznad je  $(i,j)$  element matrice  $\frac{\partial C}{\partial \theta} = \nabla C$ , a ako malo drugačije zapišemo indekse iznad, lako dolazimo do konačne vektorske formule gradijenta funkcije  $C$ :

$$11. \nabla C = \mathbf{X}^T (\mathbf{Y}' - \mathbf{Y})$$

Sa jednačinom 11, možemo kompletno da izračunamo gradijent vektor  $\nabla C = (\frac{\partial C}{\partial \theta_{0,1}}, \frac{\partial C}{\partial \theta_{1,1}}, \dots, \frac{\partial C}{\partial \theta_{k,m}})$ , i iterativno primenjivati opisanu metodu spusta gradijentom. Tehnika koju smo primenili iznad je osnovna ideja koja se krije iza algoritma propagacijom unazad (eng. *Backpropagation*). Ideja je da krenemo od zadnjeg čvora (funkcija greške), i računamo parcijalni izvod samo za taj sloj u odnosu na izlaz sloja ispod. Dobijeni rezultat propagiramo unazad.

Da bi opisani algoritam bio primenjiv u praksi, potrebno je izvršiti regularizaciju modela. Drugim rečima, potrebno je ograničiti model kako ne bi bio previše kompleksan u odnosu na trening podatke, i uveo veliku varijansu u model, tj. preveliku zavisnost od uzorka. Regularizacija prevazilazi obim ovog rada, tako da ćemo je izostaviti. Za više detalja o tehnikama regularizacije, pogledajte [12, Ch. 7]. Također, u praksi se ne vrši zbir nad čitavim skupom  $E$ , nego se iz tog skupa na slučajan način odabere manji uzorak  $B$ , i nad njim se radi spust gradijentom. Ovaj uzorak  $B$ , se ponovo uzorkuje pri svakoj iteraciji algoritma. Ova tehnika naziva se stohastički spust gradijentom [12, Ch. 8]

U ovakvom obliku, matrica  $Y'$  se dobija kao  $S(X\theta)$  gde se *softmax* primenjuje nad redovima matrice  $X\theta$ . Ovaj algoritam se naziva propagacija napred (eng. *forward pass*).

Primer koda koji implementira algoritam treniranja propagacijom unazad i algoritam evaluacije propagacijom napred:

```
def softmax(Z):
    return np.exp(Z) / np.sum(np.exp(Z), axis=1, keepdims=True)

def gradC(P, X, Y):
    n = X.shape[0]
    return (1 / n) * np.dot(X.T, (P - Y))

def L(Theta, X, Y):
    n = X.shape[0]
    prob = forward(X, Theta)
    loss = (-1 / n) * np.sum(Y * np.log(prob))
    return loss

def train(X, Y, learning_rate=0.1, iter=100):
    Theta = np.random.normal(scale=0.1, size=(X.shape[1], Y.shape[1]))

    for _ in range(iter):
        P = forward(X, Theta)
        Theta = Theta - learning_rate*gradC(P, X, Y)

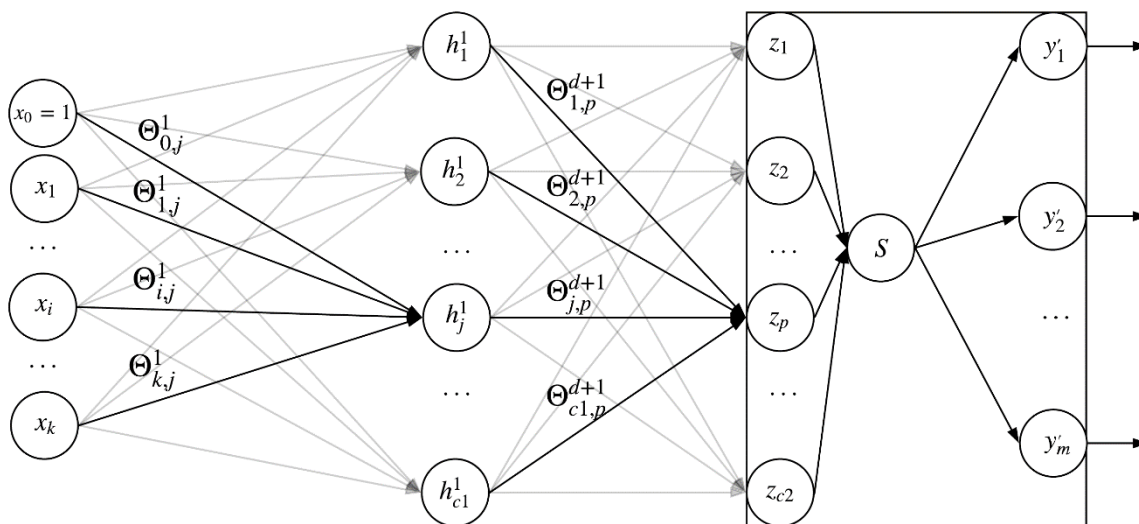
    return Theta

def forward(X, Theta):
    Z = np.dot(X, Theta)
    S = softmax(Z)
    return S
```



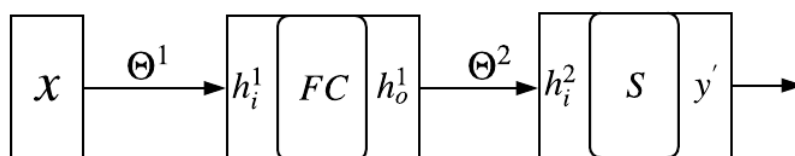
### 3.3. Duboki modeli – neuralne mreže

Neuralnu mrežu za zadatak klasifikacije dobijamo kada između čvorova ulaza i čvora *softmax* operatora ubacimo dodatne nelinearne čvorove. Ti čvorovi se nazivaju i skriveni čvorovi (eng. *hidden layers*).



Ilustracija 8: Neuralna mreža sa 1 skrivenim slojem.

Glavni deskriptivni parametar arhitekture neuralne mreže jeste dubina, tj. Broj skrivenih slojeva  $d$ . U primeru iznad, dubina je  $d=1$ . U tom slučaju, potrebne su nam 2 matrice, jedna za skriveni čvor  $\Theta^1$ , i druga za izlazni čvor  $\Theta^2$ . Sa  $c_i$  označavamo broj čvorova u  $i$ -tom skrivenom sloju.



Ilustracija 9: Kompaktni prikaz neuralne mreže. Ovakav prikaz je prikladniji za duboke neuralne mreže. Čitav sloj predstavljamo kao jedan čvor koji ima ulaz (levi odeljak), tip (srednji odeljak), i izlaz (desni odeljak). FC je skraćenica od eng. „Fully Connected“, što znači da je svaki neuron u sloju  $l$  povezan sa svakim neuronom iz sloja  $l-1$  i sloja  $l+1$ .

Glavna ideja dubokih modela jeste da svaki sloj unosi nelinearnost u model tako što se ulaz u sloj  $l$  pušta kroz aktivacionu funkciju  $f_{hl}$  koja mora da bude diferencijabilna i nelinearna funkcija. Dve funkcije koje se često koriste u praksi jesu sigmoidalna i ReLU funkcija [13], [14]. U nastavku se nećemo detaljno baviti ovim funkcijama jer ne utiču na osnovne koncepte koje ovaj rad istražuje.

Već i za mreže sa samo jednim skrivenim slojem postaje veoma teško da se izvedu jednačine za gradijent i evaluaciju direktno, pa je potreban sistematičan pristup. Umesto da računamo direktno sve slojeve, kretaćemo se jedan po jedan sloj od ulaza  $\mathbf{x}$  do izlaza prilikom evaluacije, a zatim od izlaza se vraćati ka ulazu računajući parcijalne izvode i koristeći pravilo ulančavanja izvoda složenih funkcija (ideja algoritma propagacijom unazad).

Jednačine algoritma evaluacije (eng. *forward pass*)

$$\begin{aligned} \mathbf{h}_o^1 &= f_{h1}(\mathbf{h}_i^1) \\ \text{generalno: } \mathbf{h}_o^l &= f_{hl}(\mathbf{h}_i^{l-1} \boldsymbol{\theta}^l), l = 1, \dots, d \\ \mathbf{h}_i^l &= \mathbf{h}_o^{l-1} \boldsymbol{\theta}^l \\ \mathbf{h}_o^0 &= \mathbf{x} \\ \mathbf{y}' &= S(\mathbf{h}_i^2) \end{aligned}$$

Matrica  $\boldsymbol{\theta}^l$  ima dimenzije  $c_{l-1} \times c_l$

Ako vektorizujemo algoritam sa ulaznim matricama  $\mathbf{X}$  i  $\mathbf{Y}$  kao što smo ranije opisali, algoritam evaluacije se zapisuje u sledećem rekurentnom obliku:

$$\begin{aligned} \mathbf{H}_o^l &= f_{hl}(\mathbf{H}_i^l) \text{ za } l = 1, 2, \dots, d \\ \mathbf{H}_i^l &= \mathbf{H}_o^{l-1} \boldsymbol{\theta}^l \\ \mathbf{H}^0 &= \mathbf{X} \\ \mathbf{Y}' &= S(\mathbf{H}_i^2) \end{aligned}$$

Aktivaciona funkcija  $f_{hl}$  se primenjuje na svaki element vektora/matrice posebno.

Kod ispod ilustruje algoritam evaluacije propagacijom napred:

```
def forward(X, ThetaTensor):
    H = X
    d = len(ThetaTensor) - 1

    # Compute hidden layers
    for i in range(d + 1):
        H = sigmoid(H*ThetaTensor[i]) #Sigmoid applied per element

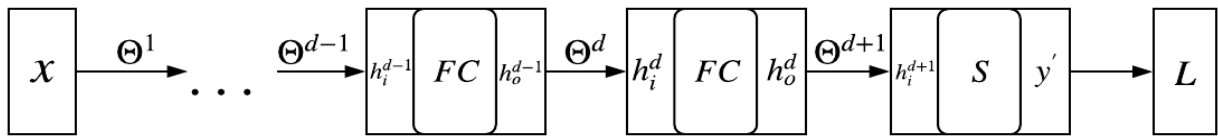
    # Compute output Softmax layer. Softmax applied per row
    return softmax(H)
```

*ThetaTensor* je niz matrica  $\boldsymbol{\theta}$  gde  $\text{ThetaTensor}[i] = \boldsymbol{\theta}^i$ . Ideja algoritma jeste da postepeno računamo vrednosti slojeva od ulaza ka izlazu, i tako propagiramo signal sve do izlaza mreže.

Za računanje gradijenta  $\nabla C$  koristićemo sličnu ideju samo što ćemo se kretati od izlaza ka ulazu. Ideja je da koristimo pravilo ulančavanja izvoda složenih funkcija kako bi kombinovali i izračunali

$$\frac{\partial C}{\partial \theta_{i,j}^l}$$

Krenućemo od izlaznog sloja, i izračunati  $\frac{\partial C}{\partial \theta^{d+1}}$ . Prilikom treniranja matrice, ubacuje se implicitan sloj koji računa unakrsnu entropiju pa treba i njega uzeti u obzir.



Ilustracija 10: Kompaktan prikaz duboke neuralne mreže. Oznake ispod strelice su izlaz iz  $i$ -tog sloja, a oznake iznad strelice predstavljaju matrice parametara.  $h_i^l$  označava ulazni vektor sloja  $l$ , a  $h_o^l$  označava izlazni vektor sloja  $l$

$$C(\theta) = \frac{1}{|E|} \sum_{(x,y) \in E} L(y, f(x; \theta))$$

gde je  $\theta$  tensor 3 reda

$$\nabla C = \frac{\partial C}{\partial \theta} = \left( \frac{\partial C}{\partial \theta^1}, \dots, \frac{\partial C}{\partial \theta^{d+1}} \right)$$

Izračunajmo:

$$\begin{aligned} \frac{\partial L}{\partial \theta^{d+1}} &= \frac{\partial L}{\partial y'} \frac{\partial y'}{\partial \theta^{d+1}} = \frac{\partial L}{\partial y'} \frac{\partial y'}{\partial h_i^{d+1}} \frac{\partial h_i^{d+1}}{\partial \theta^{d+1}} = \\ &= \frac{\partial L}{\partial h_i^{d+1}} \frac{\partial h_i^{d+1}}{\partial \theta^{d+1}} = (y' - y) \frac{\partial (h_o^d \theta^{d+1})}{\partial \theta^{d+1}} = h_o^d \otimes (y' - y) \end{aligned}$$

Ovaj izvod je identičan kao kod *softmax* klasifikatora kojeg smo izračunali ranije, pa

$$\frac{\partial L}{\partial \theta_{i,j}^{d+1}} = (y'_j - y_j) (h_o^d)_i$$

Ako označimo sa  $\delta_{d+1} = \frac{\partial L}{\partial h_i^{d+1}} = (y' - y)$  tada je:

$$\frac{\partial L}{\partial \theta^{d+1}} = h_o^d \otimes \delta_{d+1}$$

Nastavimo sa računanjem  $\frac{\partial L}{\partial \boldsymbol{\theta}^d}$

$$\frac{\partial L}{\partial \boldsymbol{\theta}^d} = \frac{\partial L}{\partial \mathbf{y}'} \frac{\partial \mathbf{y}'}{\partial \boldsymbol{\theta}^d} = \frac{\partial L}{\partial \mathbf{y}'} \frac{\partial \mathbf{y}'}{\partial \mathbf{h}_i^{d+1}} \frac{\partial \mathbf{h}_i^{d+1}}{\partial \boldsymbol{\theta}^d}$$

Izraz  $\frac{\partial L}{\partial \mathbf{y}'} \frac{\partial \mathbf{y}'}{\partial \mathbf{h}_i^{d+1}}$  smo već izračunali ranije kao  $\boldsymbol{\delta}_{d+1}$  pa ćemo ga ponovo iskoristiti (tehnika memoizacije u dinamičkom programiranju)

$$\begin{aligned} \frac{\partial L}{\partial \boldsymbol{\theta}^d} &= \boldsymbol{\delta}_{d+1} \frac{\partial \mathbf{h}_i^{d+1}}{\partial \boldsymbol{\theta}^d} = \boldsymbol{\delta}_{d+1} \frac{\partial}{\partial \boldsymbol{\theta}^d} \mathbf{h}_o^d \boldsymbol{\theta}^{d+1} = \boldsymbol{\delta}_{d+1} (\boldsymbol{\theta}^{d+1})^T \frac{\partial \mathbf{h}_o^d}{\partial \boldsymbol{\theta}^d} \\ \frac{\partial \mathbf{h}_o^d}{\partial \boldsymbol{\theta}^d} &= \frac{\partial}{\partial \boldsymbol{\theta}^d} f_{hd}(\mathbf{h}_i^d) = f'_{hd}(\mathbf{h}_i^d) \frac{\partial \mathbf{h}_i^d}{\partial \boldsymbol{\theta}^d} \end{aligned}$$

Kad uvrstimo u jednačinu iznad,

$$\frac{\partial L}{\partial \boldsymbol{\theta}^d} = \boldsymbol{\delta}_{d+1} (\boldsymbol{\theta}^{d+1})^T \circ f'_{hd}(\mathbf{h}_i^d) \frac{\partial \mathbf{h}_i^d}{\partial \boldsymbol{\theta}^d}$$

Ako pažljivije pogledamo,

$$\boldsymbol{\delta}_{d+1} (\boldsymbol{\theta}^{d+1})^T \circ f'_{hd}(\mathbf{h}_i^d) = \frac{\partial L}{\partial \mathbf{h}_i^d}$$

Označimo ovaj izraz sa

$$\boldsymbol{\delta}_d = \frac{\partial L}{\partial \mathbf{h}_i^d} = \boldsymbol{\delta}_{d+1} (\boldsymbol{\theta}^{d+1})^T \circ f'_{hd}(\mathbf{h}_i^d)$$

Jednačina  $\frac{\partial L}{\partial \boldsymbol{\theta}^d}$  postaje

$$\frac{\partial L}{\partial \boldsymbol{\theta}^d} = \boldsymbol{\delta}_d \frac{\partial \mathbf{h}_i^d}{\partial \boldsymbol{\theta}^d} = \frac{\partial}{\partial \boldsymbol{\theta}^d} \mathbf{h}_o^{d-1} \boldsymbol{\theta}^d = \mathbf{h}_o^{d-1} \otimes \boldsymbol{\delta}_d$$

Matematičkom indukcijom možemo pokazati sledeće pravilo koje nazivamo skup jednačina za propagaciju greške unazad:

$$\begin{aligned} \boldsymbol{\delta}_{d+1} &= (\mathbf{y}' - \mathbf{y}) \\ \boldsymbol{\delta}_l &= \boldsymbol{\delta}_{l+1} (\boldsymbol{\theta}^{l+1})^T \circ f'_{hl}(\mathbf{h}_i^l) \\ \frac{\partial L}{\partial \boldsymbol{\theta}^l} &= \mathbf{h}_o^{l-1} \otimes \boldsymbol{\delta}_l \end{aligned}$$

Sada kada znamo da izračunamo  $\frac{\partial L}{\partial \theta^l}$  možemo da se vratimo nazad na gradijent funkcije greške

$$\frac{\partial \mathcal{C}}{\partial \theta^l} = \frac{1}{|E|} \sum_{(x,y) \in E} \frac{\partial L}{\partial \theta^l} = \frac{1}{|E|} \sum_{(x,y) \in E} \mathbf{h}_o^{l-1} \otimes \delta_l$$

Iako ova jednačina može da se koristi za trening neuralne mreže, u praksi je jako sporo da se vrtimo petljom kroz čitav set sumirajući gradijente. Zato ćemo primeniti isti trik kao i ranije da vektorizujemo ulaz i izlaz neuralne mreže. Ako sve elemente trening skupa, spakujemo u  $|E| \times (k+1)$  matricu  $\mathbf{X}$ , a odgovarajuće koncept vektore u  $|E| \times m$  matricu  $\mathbf{Y}$ , kao što smo ranije naveli, algoritam evaluacije postaje

$$\mathbf{H}_o^l = f_{hl}(\mathbf{H}_i^l) \text{ za } l = 1, 2, \dots, d$$

$$\mathbf{H}_i^l = \mathbf{H}_o^{l-1} \theta^l$$

$$\mathbf{H}^0 = \mathbf{X}$$

$$\mathbf{Y}' = S(\mathbf{H}_i^2)$$

Ovo znači da se i vektori greške  $\delta_l$  vektorizuju u matrice i postaju  $\Delta_l$ . Trivijalnim algebarskim manipulacijama možemo jednačine propagacije unazad da vektorizujemo u sledeći oblik:

$$\Delta^{d+1} = (\mathbf{Y}' - \mathbf{Y})$$

$$\Delta^l = \Delta^{l+1} (\theta^{l+1})^T \circ f'_{hl}(\mathbf{H}_i^l)$$

$$\frac{\partial L}{\partial \theta^l} = (\mathbf{H}_o^{l-1})^T \Delta^l$$

$$\text{za } l = 1, 2, \dots, d$$

U primeru ispod je Pajton kod koji implementira potpuno povezanu neuralnu mrežu. Ovaj kod služi da bi poboljšao razumevanje jednačina iznad, i u nastavku rada se neće koristiti.

```
#Module activations.py
```

```
def softmax(Z):  
    return np.exp(Z) / np.sum(np.exp(Z), axis=1, keepdims=True)  
  
def sigmoid(Z, derivative=False):  
    S = 1/(1+np.exp(-Z))  
    return np.multiply(S, (1-S)) if derivative else S  
  
def crossentropy(Y_prim, Y):  
    return -np.sum(np.sum(Y * np.log(Y_prim), axis=1))/Y.shape[0]
```

```
#Module neural_network.py
```

```
class NeuralNetwork:  
  
    def __init__(self, c):  
        self.__c = c  
        self.__d = len(c)  
  
    @staticmethod  
    def load(path):  
        with open(path, 'rb') as handle:  
            return pk.load(handle)  
  
    def run_epoch(self, X, Y, learning_rate=0.001):  
        Theta_Grad = self.__initializeGradTensor(X.shape[1], Y.shape[1])  
        Y_prim = self.forward(X)  
  
        print("Epoch ", self.__epoch, " completed.",  
              " Loss: ", activations.crossentropy(Y_prim, Y))  
        self.__epoch = self.__epoch + 1  
  
        Delta = (Y_prim - Y)  
        Theta_Grad[self.__d+1]=np.dot(  
            activations.sigmoid(self.__H[self.__d]).T,  
            Delta)  
        self.__ThetaTensor[self.__d+1] -= learning_rate*Theta_Grad[self.__d+1]  
  
        for l in range(self.__d, 0, -1):  
            Delta = np.multiply(  
                np.dot(Delta, self.__ThetaTensor[l+1].T),  
                activations.sigmoid(self.__H[l], True))  
            Theta_Grad[l] = np.dot(activations.sigmoid(self.__H[l-1]).T, Delta)  
            self.__ThetaTensor[l] -= learning_rate*Theta_Grad[l]  
  
    def train(self, X, Y, learning_rate=0.001, epoch=100):  
        self.initialize(X.shape[1], Y.shape[1])
```

```
for e in range(epoch):
    self.run_epoch(X, Y, learning_rate)

def __initializeGradTensor(self, input_dim, output_dim):
    GradTensor = [np.zeros((1, 1))] # Dummy matrix for easier interpretation
    GradTensor.append(np.zeros((input_dim, self.__c[0])))
    for l in range(self.__d-1):
        GradTensor.append(
            np.zeros((self.__c[l], self.__c[l+1])))
    GradTensor.append(
        np.zeros((self.__c[self.__d-1], output_dim)))
    return GradTensor

def initialize(self, input_dim, output_dim):
    self.__ThetaTensor = [np.zeros((1, 1))] # Dummy matrix
    self.__ThetaTensor.append(
        np.random.normal(scale=0.1, size=(input_dim, self.__c[0])))
    for l in range(self.__d-1):
        self.__ThetaTensor.append(
            np.random.normal(scale=0.1, size=(self.__c[l], self.__c[l+1])))
    self.__ThetaTensor.append(
        np.random.normal(scale=0.1, size=(self.__c[self.__d-1], output_dim)))

    self.__epoch = 1

def save(self, path):
    with open(path, 'wb+') as handle:
        pk.dump(self, handle, protocol=pk.HIGHEST_PROTOCOL)

def forward(self, X):
    self.__H = [X]

    H_o = X
    # Compute hidden layers
    for l in range(1, self.__d+1):
        H_i = np.dot(H_o, self.__ThetaTensor[l])
        self.__H.append(H_i)
        H_o = activations.sigmoid(H_i)

    # Compute output Softmax layer
    H_i = np.dot(H_o, self.__ThetaTensor[self.__d + 1])
    self.__H.append(H_i)
    return activations.softmax(H_i)
```

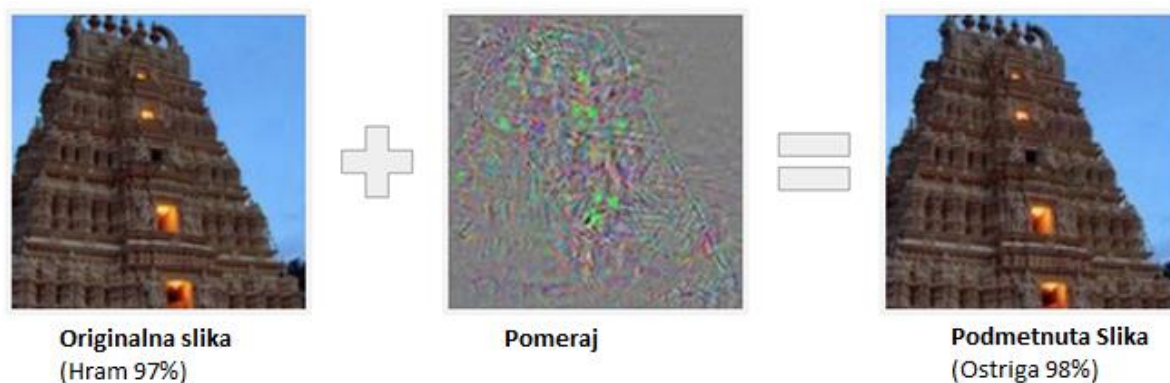
### 3.3.1. Konvolucione neuralne mreže

Neuralna mreža opisana iznad je sposobna da aproksimira bilo koju Borel-merljivu funkciju sa dovoljnim brojem skrivenih čvorova čak i sa jednim skrivenim slojem [15], međutim u praksi ove neuralne mreže sporo konvergiraju, ili se ne generalizuju dobro. Pošto ćemo se u daljem nastavku rada fokusirati na klasifikaciju slike, napomenućemo novu arhitekturu neuralne mreže koja se naziva Konvoluciona Neuralna Mreža [16].

U prvobitno opisanom modelu smo koristili isključivo potpuno povezane slojeve. To znači da je svaki neuron u sloju  $l$  povezan unapred sa svakim neuronom u sloju  $l + 1$ , itd. Konvolucione mreže uvode 2 nova sloja, konvolucioni sloj, i agregacioni sloj. Konvolucioni sloj ima zadatak da nauči filtere koji kao rezultat konvolucije pronađu korisna svojstva slike, a agregacioni sloj kombinuje pronađena svojstva sa različitih delova slike kako bi klasifikacija bila invarijantna usled translacije objekta. Uvođenje ovih slojeva ne menja opisani algoritam propagacije unazad, samo se menjaju konkretni izvodi. Pošto ovi slojevi nisu bitni za ideju, i stoga prevazilaze opseg ovog rada, izostavićemo ih. Čitalac koji želi više da sazna o konvolucionim mrežama, može da pročita na linku <http://cs231n.github.io/convolutional-networks/>, kao i u sledećim radovima [9], [12, Ch. 9], [16].

## 4. Napadi na modele dubokog učenja

Nakon što smo izložili teoretsku osnovu dubokih modela, možemo preći na glavnu ideju ovog rada, koja objašnjava kako zbuniti duboke modele da nam daju ulaz koji mi želimo. Glavna ideja kod napada koje ćemo izložiti jeste da modifikujemo ulaznu sliku u model na pametan način, tako da je model pogrešno klasifikuje, i da nam da željeni izlaz.



Ilustracija 11: Primer podmetnute slike koja kao rezultat klasifikacije daje ostrigu sa verovatnoćom 0.98. Originalna slika, i podmetnuta slika izgledaju identično za ljudsko oko.



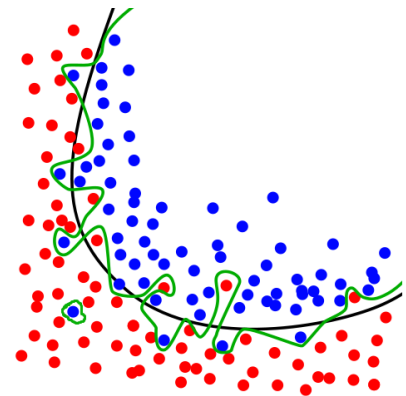
Većina modela mašinskog učenja je podložna ovim napadima, ali modeli zasnovani na gradijentnom spustu su posebno podobni za ovakve napade zbog toga što ih je veoma lako izvesti kao što ćemo kasnije videti. Također, opisani napadi se mogu izvesti na razne probleme kao što su regresija, klasifikacija, ne-nadgledano učenje (eng. *Unsupervised Learning*, učenje sa pojačavanjem (eng. *Reinforcement Learning*). Mi ćemo se u nastavku rada fokusirati na klasifikaciju slika.

Slike su po svojoj prirodi (za razliku od tekstualnog ulaza) elementi topološkog skupa (tačnije manifolda)  $I = [0, 1]^{d \times r}$  gde je  $r$  rezolucija slike, a  $d$  je dubina i uglavnom je 3 (RGB kanali). Ovo znači da ako bilo koji piksel pomerimo za neko malo  $\epsilon \in I$ , također ćemo imati validnu sliku ukoliko slika  $x + \epsilon \in I$ . Ukoliko je rezolucija slike  $r$  velika, norma vektora  $\|\epsilon\|_2$  može da bude ogromna, iako je svaki element vektora  $\epsilon$  mali.

Neka je  $\|\epsilon\|_2 = \sqrt{\sum_i \epsilon_i^2}$  L2 norma vektora  $\epsilon$ . Ako fiksiramo ovu normu na  $\|\epsilon\|_2 = p \in \mathbb{R}$ , i povećamo dimenzionalnost (rezoluciju), dopuštamo elementima vektora  $\epsilon$  da budu manji, i tako razlike između  $x$ , i  $x + \epsilon$  budu manje primetne.

Uzmimo na primer da je  $p = 1000$ . Ako je rezolucija slike 1000 piksela (ukupan broj), tada da bi prešli L2 normu od 1000 jedinica, potrebno je da se pomaknemo u proseku  $\frac{p}{1000} \sqrt{1000} \approx 31.62$  jedinica po svakoj osi. Ako je rezolucija slike 1 000 000 piksela, tada se u proseku moramo pomeriti  $\frac{p}{1000000} \sqrt{1000000} = 1$  jedinicu po svakoj osi. Ovo svojstvo može da unese probleme u model kao što ćemo ubrzo videti, i istraživači u oblasti veštačke inteligencije ga nazivaju „prokletstvo dimenzionalnosti“.

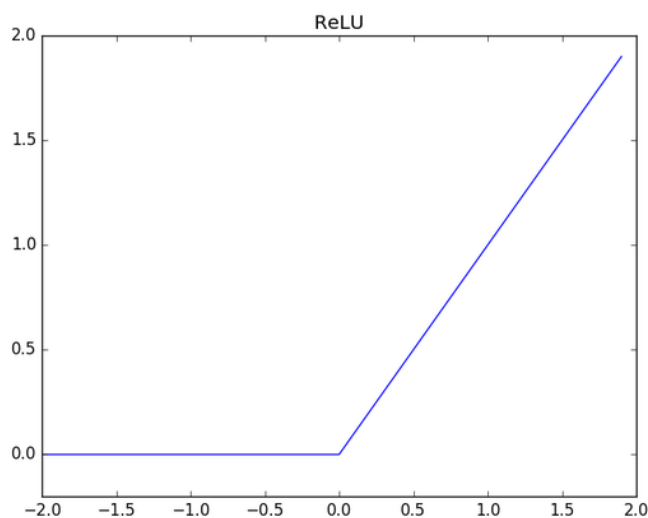
Još ranije [17], istraživači su došli na ideju da je moguće konstruisati ovakve primere koji će naterati model da proizvede željeni izlaz, ali su se uvek pitali zašto je to tako lako da se izvede kao što ćemo videti kasnije. Njihova prva hipoteza bila je da duboki modeli overfituju trening podatke (imaju veliku varijansu), i da se u prostoru slike  $I$  nalaze klasteri masa verovatnoće unutar kojih se možemo pomeriti i na taj način zbuniti model. Međutim, da je ova hipoteza tačna, tada bi svaka napadačka slika bila jedinstvena za model, ali i za svaku trening sesiju. Razlog tome je što ako model overfituje (ako mu je varijansa velika), pri svakom treniranju granica podele se znatno menja. Kao što se ispostavilo, to nije bio slučaj, i napadačke slike generisane jednom, mogu da se prenesu i da rade čak i kada se model ponovo istrenira, ali i kad se ta slika pusti kroz drugi sličan model (sa sličnom arhitekturom).



Ilustracija 12: Model je overfitovao podatke. U donjem desnom uglu vidimo jednu masu verovatnoće koja okružuje izolovani primerak plave boje.

Međutim, problem je bio što su moderni modeli linearni u delovima. Razlog ove linearnosti jeste aktivaciona funkcija koja se koristi u modernim modelima, a to je ReLU (eng. *Rectifier Linear Unit*)

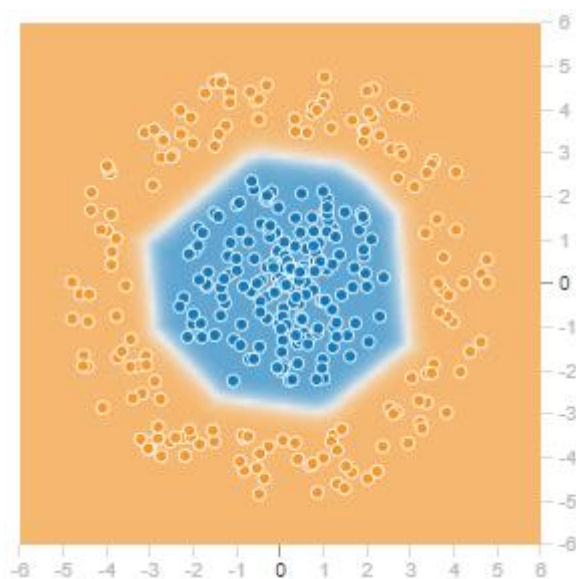
$$\text{ReLU}(x) = x^+ = \max(0, x)$$



ReLU funkcije su uvedene kako bi rešile problem gubitka gradijenta kod dubokih modela [13]. Za čitaoce koji žele više o ovome, mogu da posete link:

<https://sefiks.com/2018/05/31/an-overview-to-gradient-vanishing-problem/>

Granica podele modela istreniranih sa ReLU aktivacionom funkcijom je linearna u delovima, kao što možemo videti na slici ispod.



Ilustracija 13: Na slici vidimo da iako je granica podele nelinearna, ona je linearna u delovima

Neuralna mreža aproksimira neprekidnu granicu podele (neprekidna funkcija, ili manifold) između klasa sa konačno mnogo trening podataka. Ovo znači da neuralna mreža ekstrapolira granicu podele skoro svuda (svuda osim u konačno mnogo tačaka), i to linearno. Ova ekstrapolacija se obično dešava oko teoretske granice (funkcija koja stvarno pravi podelu između klasa, i koju ne možemo da izmerimo). Kod visoko dimenzionalnih prostora kao što je prostor slike, teoretske granice su veoma kompleksne i oblasti klasa nisu koneksne (povezane). Ako se one linearno ekstrapoliraju u tačkama gde nemamo dovoljno podataka, imamo podlogu za napadačke tehnike koje ćemo opisati kasnije.

## 4.1. Tehnike izvođenja napada

Da bi izveli napad, iskoristićemo činjenicu da je skup slika topološki prostor, i da je preslikavanje iz slike u distribuciju klasa diferencijabilno, tj.  $\mathbf{y}' = f(\mathbf{x}; \boldsymbol{\theta})$  je diferencijabilna funkcija. Ovo znači da možemo optimizovati ulaz modela  $\mathbf{x}$  tako da nam model vrati željenu vrednost vektora  $\mathbf{y}'$ .

Razlikujemo 2 vrste napada

- *ciljane napade* gde je ideja da nateramo model da nam izbaci određenu vrednost vektora  $\mathbf{y}'$
- *obične napade* gde je ideja da nateramo model da vrati pogrešnu vrednost vektora  $\mathbf{y}'$  koja god ona bila.

Napade dalje možemo grupisati u 3 grupe:

- *Analitičke napade*: Imamo potpuni pristup modelu. Tj. znamo arhitekturu modela, i imamo vrednosti svih parametara  $\boldsymbol{\theta}$ . Ova dva podatka se nazivaju i modelujući graf mreže (eng. *Computational Graph*). Ovakvi napadi su teoretski, ali služe kao osnova za ostale napade.
- *Polu analitički metod*: Nemamo potpuni pristup modelu, ali imamo pristup izlazu modela, tj. vektoru  $\mathbf{y}'$ .
- *Blackbox metod*: Nemamo pristup modelu, na izlazu dobijemo labelu kojoj klasi slika pripada.

Analitičke i polu analitičke metode izvodimo pomoću gradijenta modela zajedno sa *loss* funkcijom (slično kao pri treningu), s tim što sada ne uzimamo u obzir funkciju greške, nego direktno funkciju  $L$ .

Prilikom treninga smo definisali funkciju greške kao

$$C(\Theta) = \frac{1}{|E|} \sum_{(x,y) \in E} L(x, y; \Theta)$$

$$L(x, y; \Theta) = - \sum_{i=1}^m y_i \log f(x; \Theta)$$

Gde su  $x$ ,  $y$  bili fiksirani, a optimizacija se radila nad funkcijom  $C$  i optimizovali su se parametri  $\Theta$ . Da bi izveli napad, koristićemo sličnu ideju samo što ćemo funkciju greške definisati kao funkciju ulaza  $x$ , a parametri  $\Theta$  će biti fiksirani.

$$C(x, y) = L(x, y; \Theta) = - \sum_{i=1}^m y_i \log f(x; \Theta)$$

Da bi izveli običan napad koji nije ciljan, potrebno je da optimizujemo funkciju  $-C(x, y)$  spustom gradijenta. Dakle, za običan napad koristimo funkciju

$$J(x) = -C(x, y) = \sum_{i=1}^m y_i \log f(x; \Theta)$$

dok za ciljani napad koristimo funkciju:

$$J_t(x) = C(x, \hat{y}) = - \sum_{i=1}^m \hat{y}_i \log f(x; \Theta)$$

gde je  $\hat{y}$  željena vrednost koju želimo da nam model izbací. U nastavku ćemo se fokusirati na ciljane napade.

Optimizaciju ćemo vršiti na isti način kao i kod treniranja, tj. spustom gradijenta.

$$\begin{aligned} x'_i &= x'_i - \alpha \nabla_{x_i} J_t \\ x'_0 &= x \\ \nabla_{x_i} J_t &= \left( \frac{\partial J_t}{\partial (x'_i)_1}, \frac{\partial J_t}{\partial (x'_i)_2}, \dots, \frac{\partial J_t}{\partial (x'_i)_d} \right) \end{aligned}$$

Razlika između analitičkog i polu analitičkog metoda jeste kako računamo ove izvode. Kod analitičkog metoda računamo ih propagacijom unazad kao pri treningu, dok kod polu analitičkog metoda radimo numeričku estimaciju gradijenta. Mi ćemo se fokusirati na analitički metod, a za polu analitički metod čitalac može sam da istraži tehnike za numeričku estimaciju gradijenta.

Izračunajmo  $(\nabla_x C)_i = \frac{\partial C}{\partial (x)_i}$ . Koristeći ranije izvedene jednačine, na sličan način možemo izvesti propagaciju unazad za napad:

$$\begin{aligned}\nabla C &= \frac{\partial C}{\partial \mathbf{x}} = \boldsymbol{\delta}^1 (\boldsymbol{\theta}^1)^T \\ \boldsymbol{\delta}_{d+1} &= (\mathbf{y}' - \hat{\mathbf{y}}) \\ \boldsymbol{\delta}_l &= \boldsymbol{\delta}_{l+1} (\boldsymbol{\theta}^{l+1})^T \circ f'_{hl}(\mathbf{h}_i^l)\end{aligned}$$

$$\text{za } l = 1, 2, \dots, d$$

Za ciljane napade,  $\nabla J_t = \nabla C$ , dok je za obične napade  $\nabla J = -\nabla C$ .

Ispod je Pajton implementacija metode *run\_adversarial* za našu klasu *NeuralNetwork* koja implementira jednačine iznad:

```
def run_adversarial(self, x, y_target, num_iter=1, step_size=0.01):
    adversarial_x = x

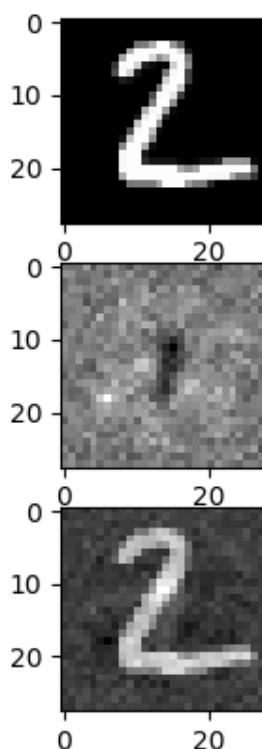
    adv_perturbation = np.zeros(adversarial_x.shape)
    for _ in range(num_iter):
        y_prim = self.forward(adversarial_x)

        delta = (y_prim - y_target)
        for l in range(self.__d, 0, -1):
            delta = np.multiply(
                np.dot(delta, self.__ThetaTensor[l+1].T),
                activations.sigmoid(self.__H[l], True))
        adv_perturbation = np.dot(delta, self.__ThetaTensor[1].T)

    adversarial_x = adversarial_x - step_size*adv_perturbation
    return adversarial_x, adv_perturbation
```

Rezultat pokretanja ovog algoritma sa sledećim parametrima je ilustrovana ispod:

```
target = [0, 1, 0, 0, 0, 0, 0, 0, 0, 0] # Digit 1  
model.run_adversarial(x, target, num_iter=7, step_size=0.1)
```



Ilustracija 14: (Gore) originalni ulaz u neuralnu mrežu klasifikovan kao cifra 2 sa verovatnoćom 0.97. (Sredina) šum koji dodajemo na ulaznu sliku. (Dole) zbir ulaza i šuma koga mreža klasifikuje kao cifru 1 sa verovatnoćom 0.77

U primeru iznad ipak vidimo razliku između ulaza i zbira ulaza i šuma. Postoje 2 razloga zašto je to tako. Prvi razlog jeste što je u primeru iznad korišćen MNIST skup za treniranje koji sadrži slike 28\*28 piksela. To nisu slike velike dimenzionalnosti pa ovaj skup i naš model sa kojim smo ga istrenirali nije podložan „prokletstvu dimenzionalnosti“ koga smo opisali ranije.

L2 norma vektora šuma  $\|\epsilon\|_2 \approx 0.928336$ , maksimalan element apsolutne vrednosti gradijenta je 0.16779047, dok je prosečan pomeraj sa ovim šumom (prosečna vrednost apsolutne vrednosti gradijenta) je 0.02517

Drugi problem jeste što naša funkcija optimizacije ne uzima u obzir koliko se slika promenila tokom optimizacije. Ovaj problem rešavamo tako što u funkciju optimizacije uvodimo regularizaciju, tj. kažemo algoritmu da ujedno pri optimizaciji kroz prostor slika  $I$

kreće tako da se što više približava našem ciljnom vektoru  $\hat{\mathbf{y}}$ , dok ujedno minimizuje  $L_p$  normu  $\|\epsilon\|_p$  slike, i tako je drži što bliže originalnoj slici. Ukoliko ulaz ima dovoljno veliku rezoluciju, ovaj metod nas vodi do rešenja koje za čoveka izgleda isto kao i originalna slika.

Reformulisaćemo našu funkciju optimizacije tako da sadrži i regularizaciju norme šuma:

$$J_{tr}(\mathbf{x}; \gamma, p) = \mathcal{C}(\mathbf{x}, \hat{\mathbf{y}}) + \gamma \|\nabla J_t\|_p$$

$$\|\mathbf{x}\|_p = \left( \sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}}$$

$$\|\mathbf{x}\|_0 = |\{ |x_i| \neq 0, i = 1, ..n \}|$$

$$\|\mathbf{x}\|_\infty = \max(|x_i|)$$

Prilikom optimizacije ove funkcije računali bi

$$\nabla J_{tr} = \nabla \mathcal{C} + \gamma \frac{\sum_i (\nabla J_{tr})_i}{\|\nabla J_{tr}\|_p} \nabla^2 J_{tr}$$

Ovaj problem minimizacije se može rešiti LBFGS algoritmom [17].

Norma  $p$  utiče na to koliko smo osetljivi na ekstreme. Što je parametar  $p$  veći to više tera algoritam da raširi pomeraje u šumu preko čitave slike, i više smo osetljivi na ekstreme. Sa normom 0 možemo da vidimo manji broj piksela koji se pomeraju, ali njihove vrednosti mogu biti velike. Sa 0 normom, možemo videti promene u slici samo u malom broju piksela ili nekoj regiji slike, dok sa  $\infty$  normom možemo zapaziti promene preko čitave slike, ali će one biti veoma male. U praksi se koristi norma  $\infty$ .

Algoritam koji ćemo sada izložiti kombinuje normu  $\infty$  i aproksimira  $\nabla J_t$ . Brza metoda znaka gradijenta koristi znak gradijenta i intenzitet pomeraja  $\epsilon$ :

$$\mathbf{x}' = \mathbf{x} - \epsilon * \text{sign}(\nabla J_t)$$

Ovaj pomeraj će da ima  $\|\mathbf{x}' - \mathbf{x}\|_\infty = \epsilon$ , i predstavlja aproksimaciju pomeraja  $\nabla J_{tr}$ . Uvođenjem momentuma [18] možemo da poboljšamo ovu metodu i da ubrzamo konvergenciju:

$$\mathbf{g}_{t+1} = \mu \mathbf{g}_t + \frac{\nabla J_t}{\|\nabla J_t\|_1}$$

$$\mathbf{x}'_{t+1} = \mathbf{x}'_t + \epsilon * \text{sign}(\mathbf{g}_{t+1})$$

$$\mathbf{g}_0 = 0$$

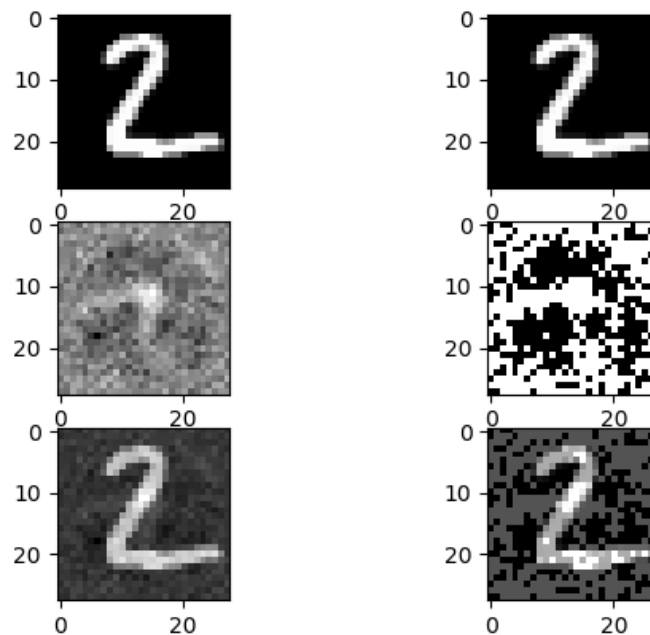
$$\mathbf{x}'_0 = \mathbf{x}$$

Ispod je primer koda koji implementira iterativni metod sa momentumom

```
def run_ifgsm(self, x, y_target, eps=1, step_size=1, momentum=0, num_iter=1):  
    r'''Iterative FGSM.  
    x: Input image  
    y_target: Optional target value  
    eps: L_{\infty} norm of the perturbation  
    num_iter: Number of iterations  
    '''  
  
    clip_min = x - eps  
    clip_max = x + eps  
  
    adversarial_x = x  
    g = 0  
    for _ in range(num_iter):  
        _, gradient = self.run_adversarial(x, y_target)  
        grad_l1norm = np.linalg.norm(gradient, 1)  
  
        g = momentum*g + (gradient/grad_l1norm)  
  
        signed_grad = np.sign(g) * step_size  
        adversarial_x = np.clip(adversarial_x + signed_grad, clip_min, clip_max)  
    adv_perturbation = adversarial_x - x  
  
    return adversarial_x, adv_perturbation
```

Iterativni metod znaka gradijenta ne radi dobro na slikama malih dimenzija jer ne može da iskoristi „prokletstvo dimenzionalnosti“ slika kao što vidimo na primeru ispod, ali se na realnim skupovima kao što su *ImageNet* i CIFAR-10 pokazuje kao bolje rešenje.



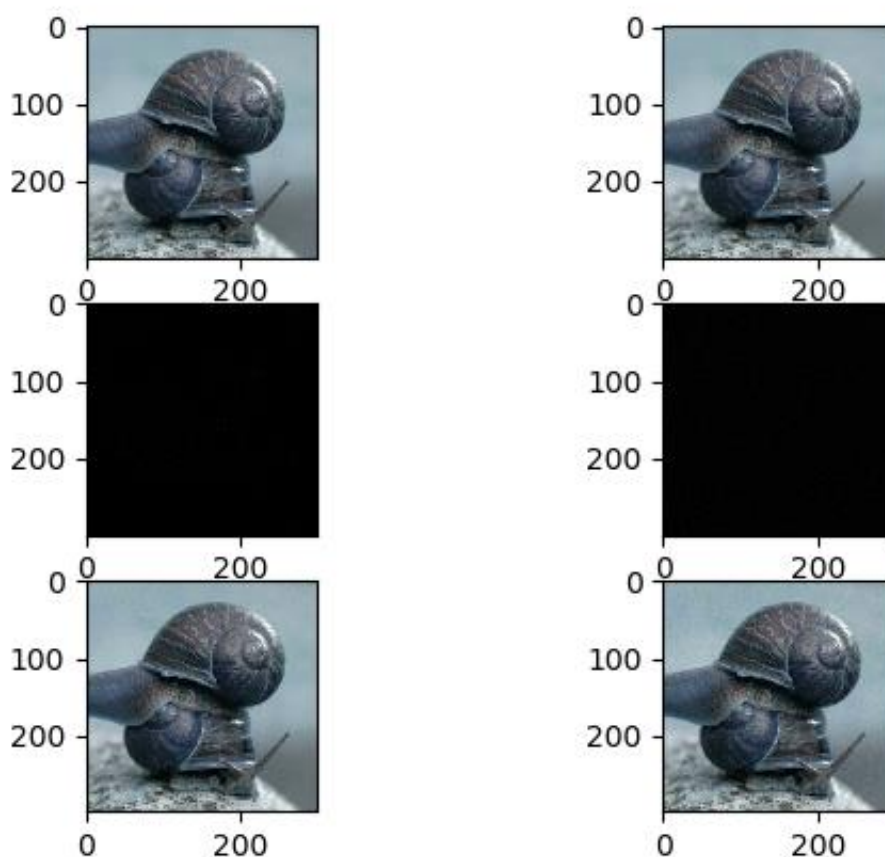


Ilustracija 15: (Leva polovina) originalna slika, šum, i na dnu rezultat spusta gradijentom. (Desna polovina) originalna slika, šum, i rezultat iterativnog metoda znaka gradijenta.

Tabela ispod prikazuje efikasnost napada na MNIST klasifikator. Originalna slika na vrhu je ulaz i model je klasifikuje kao cifru 2 sa verovatnoćom 0.97. Izlaz je modifikovana slika koja je ciljano od strane napadača da se klasifikuje kao cifra 1. Iz tabele vidimo da metod vrednosti gradijenta ima manju  $L_2$  normu od metoda znaka gradijenta, ali i veću  $L_\infty$  normu.

	Napad spustom gradijenta. Parametri: step_size(0.1) num_iter(7)	Iterativni metod znaka gradijenta. Parametri: step_size(0.1) num_iter(7) eps(0.24) momentum(0.2)
Izlaz modela nakon napada	Klasifikovano kao 1 sa verovatnoćom 0.5	Klasifikovano kao 1 sa verovatnoćom 0.52
$L_2$ norma šuma	2.20714048677303	6.719999849397163
$L_\infty$ norma šuma	0.331258678351833	0.24000000953674316
Srednja apsolutna vrednost piksela šuma	0.060804051069489815	0.23999999462132704

Na slikama veće dimenzionalnosti, model se ponaša potpuno drugačije. Ispod je rezultat primene tehnike vrednosti gradijenta kao i znaka gradijenta na *Inception v3* [19] modelu i slikama dimenzije (299, 299, 3).



Ilustracija 16: (Levo) originalna slika, šum, i rezultat na dnu posle primene metoda vrednosti gradijenta. (Desno) originalna slika, šum, i rezultat posle primene metode znaka gradijenta sa momentumom.

Tabela ispod prikazuje efikasnost napada na *ImageNet* klasifikator (*mobilenet* neuralna mreža). Originalna slika na vrhu je ulaz i model je klasifikuje kao „puž“ sa verovatnoćom 0.99. Izlaz je modifikovana slika koja je ciljana od strane napadača da se klasifikuje kao zmaj „kite“.

Možemo videti da je u proseku srednja vrednost šuma kod metoda vrednosti gradijenta manja, ali je maksimalna vrednost veća. Ovo je očekivano ponašanje pošto metoda znaka gradijenta optimizuje i vrši regularizaciju u  $L_\infty$  normi, pa samim tim model raspoređuje pomeraje šuma na sve piksele. Također možemo primetiti kako za sličnu vrednost  $L_2$  norme, ostali parametri su znatno manji nego za primere iznad sa MNIST podacima. Razlog tome jeste dimenzionalnost slika.

Originalna slika klasifikovana kao puž sa verovatnoćom 0.99	Napad spustom gradijenta. Parametri: step_size(0.1) num_iter(5)	Iterativni metod znaka gradijenta. Parametri: step_size(0.05) num_iter(5) eps(0.016) momentum(0.8)
Izlaz modela nakon napada	Klasifikovano kao zmaj („kite“) sa verovatnoćom 0.99	Klasifikovano kao zmaj („kite“) sa verovatnoćom 0.99
$L_2$ norma šuma	4.2268567	8.286149
$L_\infty$ norma šuma	0.085820615	0.016000003
Srednja apsolutna vrednost piksela šuma	0.0056709372	0.015999993

## 4.2. Tehnike odbrane

Odbrana od napada opisanih u ovom radu još uvek je otvoreno istraživačko pitanje, i sve poznate metode ne garantuju sigurnost. Ipak, metod koji je najkorisniji jeste suparničko treniranje. Ideja je da istreniramo model nad nekim skupom, zatim nad tim skupom generišemo protivničke slike kao što smo ranije opisali i te slike ubacimo u novi skup za trening. Na kraju ponovo istreniramo model sa novim skupom. Iako ova tehnika ne unosi sigurnost u model, veoma je dobra za generalizaciju modela (smanjenje njegove varijanse). Napadač koji ima više procesorske moći može jednostavno da ponovi sve navedene tehnike više puta i napadne model sa jačim napadima. Pored suparničkog treniranja postoje i razne druge tehnike koje su takođe neuspešne [20], [21].

Postavlja se pitanje zašto je odbrana od ovih modela toliko teška. Naime, od modela se zahteva da daju dobar izlaz na svaki mogući ulaz, ali se za njihovo treniranje koristi konačno mnogo podataka. Dalje, probabilistički modeli uče da daju izlaz tako što aproksimiraju distribuciju podataka koje su videli tokom treninga, dakle, distribuciju stvarnih slika koje se javljaju u prirodi. Kada modelu prikažemo sliku sa dodatnim šumom koji nije beli šum, menja se distribucija i narušavaju se pretpostavke koje je model uspostavio tokom treninga. Drugi problem je nedostatak teoretskih modela koji bi opisali šta se tačno dešava nakon optimizacije. Ovakvi modeli su neophodni da bi mogli da konstruišemo argument i

odgovorimo na pitanje „šta je suparnički primer?“. Bez ovakvih modela nismo u stanju da utvrdimo da li bilo koja tehnika odbrane daje rešenje ili ne.

Primer koji najbolje ilustruje zašto se ovakvi napadi mogu lako konstruisati jeste primer „Pametnog Hansa“. Pametni Hans je bio konj u doba renesanse koji je navodno znao osnovnu aritmetiku, tj. da sabira brojeve. Na glavnom trgu grada bi se okupili svi građani i zadavali konju aritmetičke zadatke. Jedan građanin bi rekao „Hans, koliko je  $2+3$ ?“, na šta bi Hans udario kopitom od tlo 5 puta i u tom trenutku bi publika počela da aplaudira.

Jedan psiholog nije mogao da veruje šta se dešava pa je uzeo Hansa da ga testira. Kada ga je doveo u svoje prostorije, i izolovao od svega, Hans nije više mogao da radi aritmetiku. Ispostavilo se da Hans nije znao da radi aritmetiku nego je posle svakog pitanja tapkao dok ne čuje aplauze publike. Ni Hans ni njegov pastor nisu bili prevaranti, nego je Hans naučio da rešava problem u okruženju u kome je treniran i u kome je stalno testiran. To okruženje je bilo ono gde je Hans bio okružen publikom koja bi počela da aplaudira kada bi Hans udario određen broj puta kopitom o tlo. Kada je Hans izolovan iz takvog okruženja, nije više mogao da obavlja zadatak.

## 5. Zaključak

Duboki modeli leže u srcu moderne veštačke inteligencije, i možemo ih sresti skoro u bilo kom modernom proizvodu. Nalaze se u pametnim telefonima kako bi poboljšali kameru, i prepoznali vlasnika, u računarima kako bi prepoznali govor. Također, razni servisi na internetu koriste duboke modele kako bi pružili što priyatnije iskustvo za korisnika. Odgovarajuće mere sigurnosti su potrebne kako korisnik modela ne bi pretrpeo štetu nastalu usled malicioznog napada. Kao što smo ranije pokazali, ove napade je veoma lako izvesti, pa čak i onda kada nemamo direktan pristup modelu tako što numerički procenimo gradijente. Ipak prava opasnost leži od toga što je moguće konstruisati primer za jedan model (neki od popularnih modela koje možemo da treniramo kući), i preneti ga na sličan model koji radi u pozadini servisa [22].

Sve navedene ranjivosti dubokih modela treba da nam ukažu da je sigurnost veštačke inteligencije ozbiljan i nerešen problem, te je potrebno ozbiljno se pozabaviti istim pre nego što čovečanstvo napravi sledeći korak prema opštoj veštačkoj inteligenciji. Iako trenutno ne postoji efikasna odbrana od tehnika napada izloženih u ovom radu, ipak nije sve izgubljeno. Napadi se znatno mogu otežati ako se onemogućiti napadaču da vidi verovatnoće koje model izbacuje, kao i ako se ograniči broj upita ka servisu koji pruža usluge veštačke inteligencije. Ove metode u ovom momentu mogu da otežaju znatno napade, ali ni one nisu sasvim otporne. Također je bitno napomenuti da su probabilistički modeli za klasifikaciju slika najviše ranjivi jer su slike topološki prostori sa velikom dimenzionalnošću. Dakle, redukcija rezolucije može također da se koristi kao tehnika odbrane.

Od tehnika koje smo naveli, metoda znaka gradijenta daje bolje rezultate za fiksnu jedinicu vremena nego metoda vrednosti gradijenta jer brže konvergira i uzima u obzir i regularizaciju. Obe metode mogu da budu moćno oružje u rukama napadača, i mogu da naštetu sistemima koji su kritični kao što su automobili, industrija, i slično. Trenutno su ovi modeli dobro zaštićeni i još uvek je teže napasti takve sisteme u praksi, ali sa vremenom će to postajati lakše. Iako ovaj problem nije prioritetan, da bi uveli veštačku inteligenciju u svakodnevni život, i prepustili joj da obavlja kritične zadatke, pitanje odbrane od ovakvih napada sigurno treba da se odgovori.

## 6. Reference

- [1] S. Hoehl, K. Hellmer, M. Johansson, and G. Gredebäck, "Itsy Bitsy Spider...: Infants React with Increased Arousal to Spiders and Snakes," *Front. Psychol.*, vol. 8, Oct. 2017.
- [2] F. T. Bruss, "250 years of 'An essay towards solving a problem in the doctrine of chances. By the late Rev. Mr. Bayes, F.R.S. communicated by Mr. Price, in a letter to John Canton, A.M.F.R.S.,'" 2014.
- [3] W. S. McCulloch and W. Pitts, "A logical calculus of the ideas immanent in nervous activity," *Bull. Math. Biophys.*, vol. 5, no. 4, pp. 115–133, Dec. 1943.
- [4] A. M. TURING, "I.—COMPUTING MACHINERY AND INTELLIGENCE," *Mind*, vol. LIX, no. 236, pp. 433–460, Oct. 1950.
- [5] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *Nature*, vol. 323, no. 6088, pp. 533–536, Oct. 1986.
- [6] J. J. Hopfield, "Neural networks and physical systems with emergent collective computational abilities.," *Proc. Natl. Acad. Sci. U. S. A.*, vol. 79, no. 8, pp. 2554–8, Apr. 1982.
- [7] Y. Yann, Y. Yann, L. Le, C. Cun, and S. Trainable, "Yann LeCunThe The Challenges of of Machine Learning."
- [8] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997.
- [9] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks."
- [10] R. DiPietro, "A Friendly Introduction to Cross-Entropy Loss." [Online]. Available: <https://rdipietro.github.io/friendly-intro-to-cross-entropy-loss/>. [Accessed: 23-Sep-2018].
- [11] D. J. C. MacKay, *Information theory, inference, and learning algorithms*. Cambridge University Press, 2003.
- [12] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. .
- [13] P. Ramachandran, B. Zoph, and Q. V. Le, "Searching for Activation Functions," Oct. 2017.
- [14] K. He, X. Zhang, S. Ren, and J. Sun, "Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification," Feb. 2015.
- [15] K. Hornik, "Approximation capabilities of multilayer feedforward networks," *Neural Networks*, vol. 4, no. 2, pp. 251–257, Jan. 1991.
- [16] Y. LeCunn and Y. Bengio, "Convolutional Networks for Images, Speech, and Time'Series," 1995.
- [17] C. Szegedy *et al.*, "Intriguing properties of neural networks."
- [18] Y. Dong *et al.*, "Boosting Adversarial Attacks with Momentum," Oct. 2017.
- [19] Z. Alom, M. Hasan, C. Yakopcic, T. M. Taha, and V. K. Asari, "Improved Inception-Residual Convolutional Neural Network for Object Recognition."

- [20] N. Papernot, P. McDaniel, X. Wu, S. Jha, and A. Swami, “Distillation as a Defense to Adversarial Perturbations against Deep Neural Networks,” Nov. 2015.
- [21] N. Papernot, P. McDaniel, A. Sinha, and M. Wellman, “Towards the Science of Security and Privacy in Machine Learning,” Nov. 2016.
- [22] F. Tramèr, N. Papernot, I. Goodfellow, D. Boneh, and P. McDaniel, “The Space of Transferable Adversarial Examples,” Apr. 2017.