

SVEUČILIŠTE U ZAGREBU
FAKULTET ELEKTROTEHNIKE I RAČUNARSTVA

SEMINAR

Primjena transformer modela za klasifikaciju slika

Rej Šafranko

Voditelj: *prof. dr. sc. Siniša Šegvić*

Zagreb, svibanj 2023.

SADRŽAJ

1. Uvod	1
2. Model zasnovan na pažnji	2
2.1. Koder-dekoder arhitektura	3
2.2. Mehanizam pažnje	3
2.3. Pozicijsko kodiranje	4
3. Model za vid zasnovan na pažnji	6
3.1. Ulaz u model	7
3.2. Koder modela	7
3.3. Induktivna pristranost modela	8
4. Certificirana robustnost klasifikacije slika	9
4.1. Ablacije slike	9
4.2. Derandomizirano zaglađivanje	10
4.2.1. Zaglađeni klasifikator	10
4.2.2. Certificirana robusnost zaglađenog klasifikatora	10
4.3. Zaglađeni model za vid zasnovan na pažnji	11
5. Reproduciranje rezultata	12
5.1. Programska implementacija	12
5.2. Skup podataka CIFAR-10	13
5.3. Rezultati	13
6. Klasifikacija vrsta riba	15
6.1. Skup podataka	15
6.2. Korištene tehnologije i parametri učenja	15
6.3. Rezultati	16

7. Zaključak	17
8. Literatura	18
9. Sažetak	19

1. Uvod

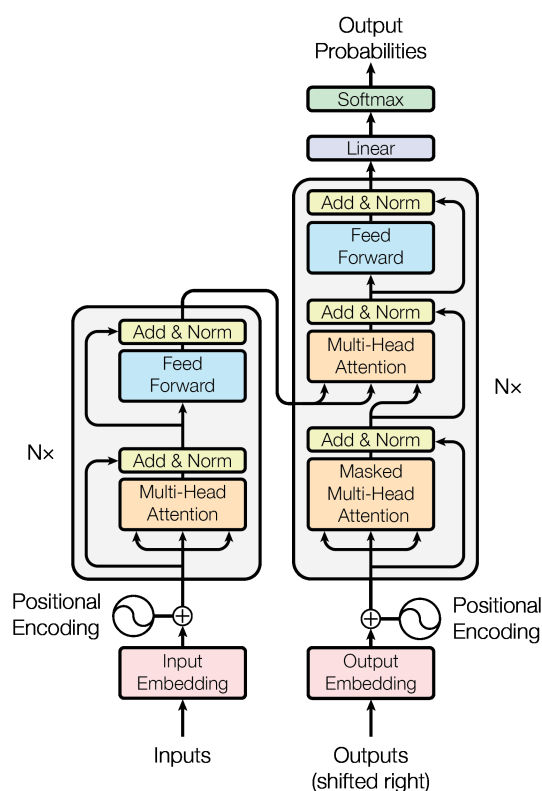
Vrhunski rezultati u obradi prirodnog jezika ostvaruju su uporabom modela zasnovanim na pažnji. Postavlja se pitanje mogu li se modeli zasnovani na pažnji značajno primjeniti u računalnom vidu. Preciznije, je li moguće ostvariti rezultate u zadatku klasifikacije slika koji se mogu mjeriti sa rezultatima konvolucijskih neuronskih mreža, koje dominiraju područjem računalnog vida (1)?

Primjena računalnog vida u visokorizičnim situacijama zahtjeva razvoj sustava koji su garantirano robusni na nepredvidive promjene (ablacije) u ulaznim podacima. Sustavi računalnog vida griješe u klasifikaciji kad su izloženi kontradiktornim napadima (staklo, grafiti, odjeća, naljepnice). Obranu od kontradiktornih napada je teško ocijeniti jer se napade može prilagoditi i zakomplicirati. Ovo dovodi do potrebe za sustavima koji su robusni na napade bez empirijske evaluacije. Uporaba modela za vid zasnovanim na pažnji je jedno od mogućih rješenja ovog problema (2).

U ovom seminarskom radu istražiti ću arhitekturu modela zasnovanim na pažnji, prilagodbu modela za zadatke računalnog vida te konkretan problem ostvarivanja garantirane robusne klasifikacije slika uporabom modela za vid zasnovanom na pažnji. Reproducirat ću rezultate eskperimenta iz (2) na skupu podataka CIFAR-10 (3). Na kraju ću istražiti učinkovitost modela za vid zasnovanom na pažnji u problemu klasifikacije vrsta riba.

2. Model zasnovan na pažnji

Modeli zasnovani na pažnji su trenutno dominantni modeli u području obrade prirodnog jezika koji modeliraju sekvence. Prije se taj zadatak obavljao konvolucijskim i povratnim modelima izvedenim kao koder-dekoder arhitekture. Modeli zasnovani na pažnji se također zasnivaju na koder-dekoder arhitekturi, ali ne koriste ni konvolucijske slojeve ni povratne veze za modeliranje sekvenci. Koriste mehanizam pažnje za modeliranje ovisnosti dijelova sekvenci, neovisno o njihovoj udaljenosti unutar sekvence (globalno receptivno polje) (4).



Slika 2.1: Arhitektura modela zasnovanog na pažnji (4)

2.1. Koder-dekoder arhitektura

Koder modela zasnovanog na pažnji se sastoji od 6 identičnih slojeva, a svaki sloj ima 2 podsloja. Prvi podsloj čini mehanizam pažnje s više glava, a drugi čini potpuno povezana unaprijedna neuronska mreža. Svaki podsloj ima rezidualnu konekciju (5) koja, zajedno sa izlazom tog podsloja, ulazi u normalizacijski sloj (6).

Dekoder modela zasnovanog na pažnji je građen kao i koder, no sadrži treći podsloj koji izvodi mehanizam pažnje s više glava nad izlazom koder modela. Koriste se rezidualne konekcije za svaki podsloj na isti način kao i kod koder modela. Mehanizam pažnje kod dekoder modela je prilagođen tako da na izlaz modela za neku poziciju sekvence mogu utjecati samo izlazi prijašnjih pozicija te sekvence (4).

2.2. Mehanizam pažnje

Mehanizam pažnje preslikava nizove vektora upita Q , ključa K i vrijednosti V u izlazni niz vektora.

U **dekoderu** vektor Q dolazi od prijašnjih slojeva dekoder modela, a vektori K i V dolaze od izlaza koder modela. Ovo omogućuje da pozicija sekvence u dekoderu "obraća pažnju" na ostale (prijašnje) pozicije u sekvenci (4).

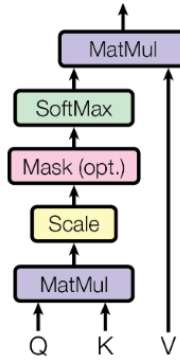
U **koderu** vektori Q , K i V dolaze iz prijašnjeg sloja koder modela. Pozicija sekvence u koderu "obraća pažnju" na pozicije sekvence u svim prijašnjim slojevima koder modela (4).

Ulaz u funkciju pažnje čine vektori Q i K dimenzije d_k te vektor V dimenzije d_v . Potrebno je izračunati softmax skalarnog produkta između Q i K podijeljenog sa $\sqrt{d_k}$. Na kraju se izračunaju težine kao skalarni produkt izlaza softmaxa i vektora V .

$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}})V \quad (2.1)$$

Ovaj postupak se može provesti više puta. Vektori Q , K i V dimenzije d_{model} se preslikavaju (sa različitim i naučenim linearnim preslikavanjima) h puta na dimenzije d_k , d_k i d_v . Sada se nad izlazom svakog preslikavanja paralelno provodi funkcija pažnje koja daje izlazne vektore (glave) dimenzije d_v . Glave se konkatenuiraju i ponovno preslikavaju što rezultira konačnim izlazom.

$$MultiHead(Q, K, V) = Concat(head_1, \dots, head_h)W^0 \quad (2.2)$$



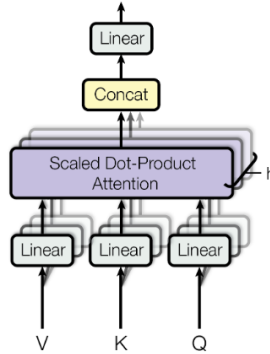
Slika 2.2: Mehanizam pažnje (4)

gdje je

$$head_i = Attention(QW_i^Q, KW_i^K, VW_i^V) \quad (2.3)$$

a preslikavanja su matrice parametara $W_i^Q \in \mathbb{R}^{d_{model} \times d_k}$, $W_i^K \in \mathbb{R}^{d_{model} \times d_k}$ i $W_i^V \in \mathbb{R}^{d_{model} \times d_v}$ (4).

Opisani postupak se naziva mehanizam pažnje s više glava.



Slika 2.3: Mehanizam pažnje s više glava (4)

2.3. Pozicijsko kodiranje

Pošto model zasnovan na pažnji ne sadrži ni povratne veze ni konvolucijske slojeve, potrebno je unjeti informacije o relativnoj ili apsolutnoj udaljenosti pozicija u sekvenci kako bi model mogao iskoristiti poredak tokena u sekvenci (4). Prije ulaska u koder ili dekodek, ulazna kodiranja se sumiraju sa pozicijskim kodiranjima. Ulazna i pozicijska kodiranja su jednake dimenzije d_{model} .

Postoje razna pozicijska kodiranja koja mogu biti ili naučena ili fiksna (7). U radu (4) koji predstavlja transformer model, korištena su trigonometrijska kodiranja različitih frekvencija. Pretpostavka je da će odabirom ovih funkcija model lakše naučiti relativne odnose pozicija u sekvenci (4).

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{model}}) \quad (2.4)$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{model}}) \quad (2.5)$$

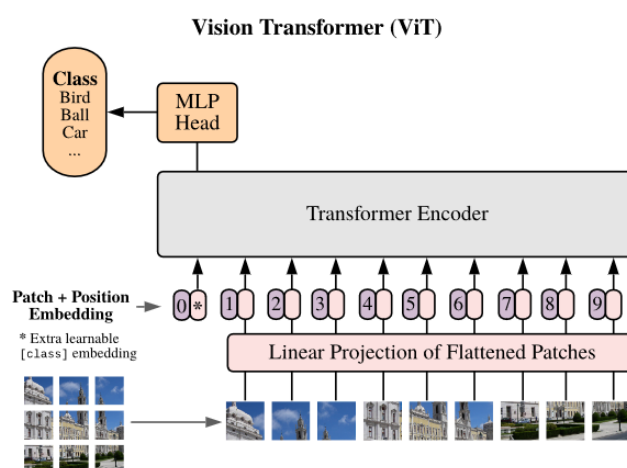
gdje pos predstavlja poziciju u sekvenci, a i dimenziju pozicijskog kodiranja.

3. Model za vid zasnovan na pažnji

U području računalnog vida dominantni modeli su konvolucijske neuronske mreže (1). Nakon uspjeha modela zasnovanim na pažnji u obradi prirodnog jezika, nekoliko radova je pokušalo povezati konvolucijske arhitekture i mehanizam pažnje (8; 9). To je inspiriralo razvoj modela za vid zasnovanim na pažnji koji ima minimalne promjene u odnosu na standardni model zasnovanim na pažnji. (1).

Arhitektura modela za vid zasnovanim na pažnji je u skladu s arhitekturom standardnog modela zasnovanim na pažnji. Najveća razlika se očekuje na ulazu modela pošto model radi sa slikama, a minimalna razlika u koder-dekoder arhitekturi.

Slika se podijeli u isječke fiksne veličine. Isječci se linearno preslikaju u kodiranje. Kodirani isječci se sumiraju sa pozicijskim kodiranjima. Dobivena sekvenca vektora čini ulaz u model. Da bi radili klasifikaciju, potrebno je dodati klasifikacijski token u sekvencu.



Slika 3.1: Model za vid zasnovan na pažnji (1)

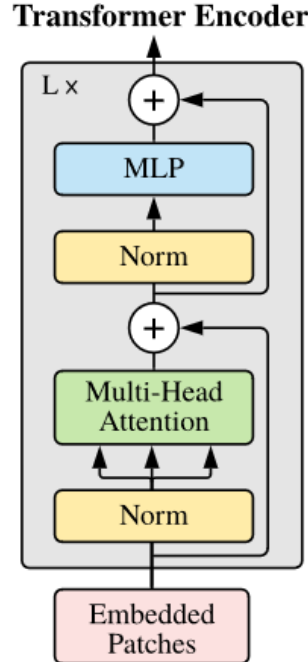
3.1. Ulaz u model

Standardni model zasnovan na pažnji prima ulaz kao 1D sekvencu tokena (ulazno kodiranje). Kako bi model mogao raditi sa slikama (2D ulazima), potrebno je promijeniti dimenzije slike. Slika $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$ se pretvori u sekvencu $\mathbf{x}_p \in \mathbb{R}^{N \times (P^2 \cdot C)}$ koja sadrži 2D isječke slike. (H, W) predstavlja rezoluciju originalne slike, a C je broj kanala. (P, P) je rezolucija svakog 2D isječka originalne slike, a $N = HW/P^2$ je broj isječaka. N je ujedno i duljina ulazne sekvence.

Model radi sa vektorima dimenzije D kroz sve slojeve pa se isječki izravnavaju i mapiraju na na D dimenzija sa naučenim linearnim preslikavanjem. Izlazi preslikavanja zovu se kodirani isječki.

Pozicijska kodiranja sumiraju se sa kodiranim isječcima radi održavanja prostorne informacije. Također, u sekvencu se ubacuje klasifikacijski token x_{class} koji se može naučiti (1).

3.2. Koder modela



Slika 3.2: Koder modela za vid zasnovanim na pažnji (1)

Koder modela sadrži alternirajuće slojeve mehanizma pažnje s više glava i MLP blokova. Normalizacija sloja se primjenjuje prije svakog bloka, a rezidualne konekcije se

dodaju izlazu svakog bloka (1).

3.3. Induktivna pristranost modela

Bitno je napomenuti da je induktivna pristranost specifična za slike manja kod modela za vid zasnovanim na pažnji nego kod konvolucijskih neuronskih mreža. Kod konvolucijskih mreža lokalnost i struktura 2D susjedstva su prisutni kroz sve slojeve modela. Kod modela za vid zasnovanim na pažnji, lokalnost je prisutna samo kod MLP slojeva, dok su slojevi mehanizma pažnje globalni. Struktura 2D susjedstva je prisutna samo na početku kod podjele slike na isječke. Inicijalna pozicijska kodiranja ne nose nikakvu informaciju o relativnim i apsolutnim pozicijama isječaka u 2D susjedstvu, već se ona moraju naučiti (1).

Jedno alternativno formiranje ulazne sekvence je korištenje mapi značajki konvolucijskih neuronskih mreža umjesto isječaka originalne slike. Ovakav hibridni model bi mogao imati induktivnu pristranost više specifičnu za slike.

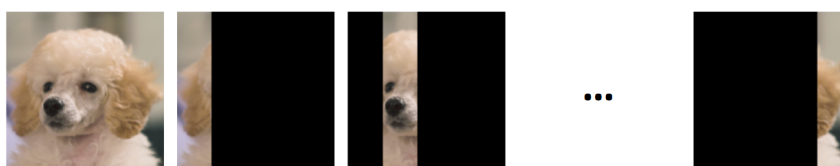
4. Certificirana robustnost klasifikacije slika

Sustavi računalnog vida koriste se u situacijama koje su visokog sigurnosnog rizika (npr. cestovni i zračni promet) te moraju biti pouzdani u svakom mogućem scenariju (2). Jedan od glavnih problema u računalnom vidu su smetnje u podacima. To može biti neka klasa koju model još nije vidio (izvandistribucijska klasa) ili neprijateljski primjeri. Neprijateljski primjeri su skup pojava koje zavaraju sustav računalnog vida (staklo, grafiti, naljepnice, odjeća) te sustav čini pogrešku u klasifikaciji. Jednostavan primjer neprijateljskog primjera je čovjek na cesti u odjeći koja se uklapa sa okolinom ceste.

Aktivna tema istraživanja u računalnom vidu je razvoj sustava koji su garantirano otporni, tj. robusni, na neprijateljske primjere. Za takve sustave se kaže da su certificirano robusni. Jedan način postizanja certificirano robusne klasifikacije je korištenje modela za vid zasnovanima na pažnji kao osnovicu za obranu zaglađivanjem. Promatrana obrana zaglađivanjem se zove derandomizirano zaglađivanje i temelji se na predikcijama klasifikatora na ablacijama ulazne slike (2).

4.1. Ablacije slike

Ablacije slike su varijacije slike kod kojih je većinski dio slike maskiran. Postoje stupčaste ablacije i ablacije u blokovima. Dalje razmatram samo stupčaste ablacije.



Slika 4.1: Stupčaste ablacije (2)

Kod stupčastih ablacija, nemaskirani dio slike je fiksne veličine. Za sliku \mathbf{x} rezolucije $h \times w$, skup svih mogućih stupčastih ablacija širine b je $S_b(\mathbf{x})$. Stupčasta ablacija može početi na bilo kojoj poziciji od njih w i može biti omotana oko slike što znači da je w ukupni broj ablacija u $S_b(\mathbf{x})$ (2).

4.2. Derandomizirano zaglađivanje

Derandomizirano zaglađivanje je jedan od načina obrane zaglađivanjem koji konstruira zaglađeni klasifikator. Taj klasifikator se sastoji od dva glavna dijela: baznog klasifikatora i skupa ablacija slike $S_b(\mathbf{x})$. Bazni klasifikator se zaglađuje tim skupom ablacija. Rezultirajući zaglađeni klasifikator vraća najčešću predikciju baznog klasifikatora na skupu ablacija $S_b(\mathbf{x})$ (2).

4.2.1. Zaglađeni klasifikator

Za ulaznu sliku x , skup ablacija $S_b(\mathbf{x})$ te bazni klasifikator f , zaglađeni klasifikator g se definira kao:

$$g(\mathbf{x}) = \operatorname{argmax}_c n_c(\mathbf{x}) \quad (4.1)$$

gdje je

$$n_c(\mathbf{x}) = \sum_{\mathbf{x}' \in S_b(\mathbf{x})} [\mathbb{I}f(\mathbf{x}') = c] \quad (4.2)$$

broj ablacija slike koje su klasificirane kao klasa c . \mathbb{I} predstavlja indikatorsku funkciju. Udio slika koje je zaglađeni klasifikator točno klasificirao zove se standardna preciznost (2).

4.2.2. Certificirana robusnost zaglađenog klasifikatora

Zaglađeni klasifikator je certificirano robusan za ulaznu sliku ako je broj ablacija klasificiranih u najčešću klasu veći od broja ablacija klasificiranih u drugu najčešću klasu za dovoljno veliku marginu. Velika margina onemogućava neprijateljskom primjeru da promjeni predikciju zaglađenog klasifikatora jer jedan neprijateljski primjer može biti pristuan samo na ograničenom broju ablacija (2).

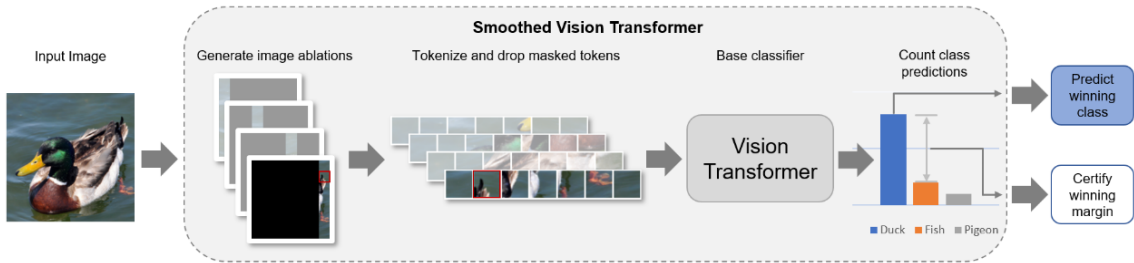
Formalno, neka je Δ maksimalni broj ablacija u skupu ablacija $S_b(\mathbf{x})$ koje jedan neprijateljski primjer može obuhvatiti istovremeno (za stupčaste ablacije širine b , neprijateljski primjer rezolucije $m \times n$ može obuhvatiti najviše $\Delta = m + b - 1$ ablacija).

Onda je zaglađeni klasifikator certificirano robusan na ulazu x ako za predviđenu klasu c vrijedi:

$$n_c(\mathbf{x}) > \max_{c' \neq c} n_{c'}(\mathbf{x}) + 2\Delta \quad (4.3)$$

Ako se zadovolji ovaj uvjet, najčešća klasa će garantirano ostati nepromijenjiva čak i ako neprijateljski primjer kompromitira svaku ablaciju koju obuhvati (2).

4.3. Zaglađeni model za vid zasnovan na pažnji



Slika 4.2: Zaglađeni model za vid zasnovan na pažnji (2)

Model za vid zasnovan na pažnji posjeduje 2 svojstva koja ga čine pogodnim za derandomizirano zaglađivanje:

1. Obraduje sliku kao sekvencu isječaka te slike. Iz toga slijedi da može odbaciti nepotrebne isječke iz sekvence i ignorirati veće dijelove slike.
2. Mehanizam pažnje u modelu dijeli informacije globalno kroz sve slojeve. To je pogodno za klasifikaciju ablacija jer će model pridati više pažnje manjem nemaskiranom dijelu slike.

Korištenjem modela za vid zasnovanim na pažnji kao bazni klasifikator te zaglađivanjem istog ablacijama, dobiva se zaglađeni model za vid zasnovan na pažnji.

Korisno je povući poveznicu sa konvolucijskim neuronskim mrežama kako bi se opravdala uporaba modela za vid zasnovanim na pažnji. Razlika je što konvolucijske mreže ne obrađuju sliku kao sekvencu isječaka pa se receptivno polje, koje je lokalno, mora postepeno izgraditi. To povlači da konvolucijske mreže moraju obraditi i maskirane dijelove slike koje model za vid zasnovan na pažnji ignorira. Sada je intuitivno jasnije da je model za vid zasnovan na pažnji pogodniji odabir za postizanje certificirane robusnosti klasifikacije slika.

5. Reproduciranje rezultata

U ovom poglavlju prikazujem ishod reprodukcije rezultata rada (2) koji predstavlja zaglađeni model za vid zasnovan na pažnji za postizanje certificirane robusnost klasifikacije.

5.1. Programska implementacija

Kako bih reproducirao rezultate rada (2), prilagodio sam originalni Python kod autora kao Jupyter bilježnicu. Korišteni radni okviri su Pytorch i MadryLab Robustness. Jupyter bilježnicu pokrećem na platformi Google Colaboratory koja nudi besplatan pristup grafičkoj procesnoj jedinici NVIDIA Tesla T4.

Učim ViT-T (2) model na skupu podataka CIFAR-10 (3). Parametri treniranja su sljedeći:

- Broj epoha: 30
- Optimizator: AdamW
- Stopa učenja: 0.01
- Propadanje težina: $5 \cdot 10^{-4}$
- Veličina grupe: 128
- Vrsta ablacija: stupčasta
- Širina ablacija (b): 4

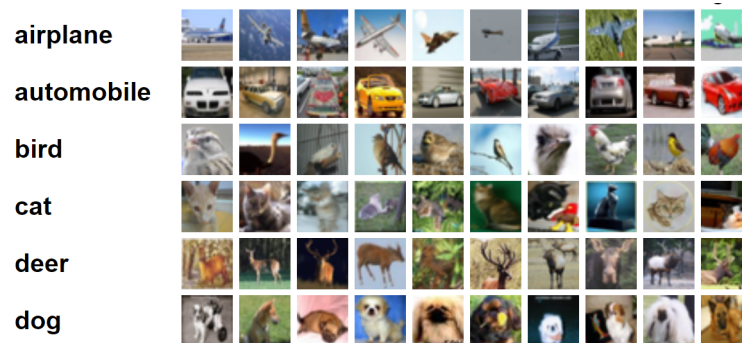
Nakon treniranja treba certificirati robusnost modela. Parametri certifikacije su:

- Vrsta ablacija: stupčasta
- Širina ablacija (b): 4
- Veličina neprijateljskih primjera (m): 5

5.2. Skup podataka CIFAR-10

CIFAR-10 sadrži 60 000 RGB slika rezolucije 32×32 . Slike su podijeljene u 10 klasa te svaka klasa ima 6000 slika. CIFAR-10 se dijeli na skup za učenje od 50 000 slika i skup za ispitivanje od 10 000 slika. Skup za ispitivanje se sastoji od po 1000 slučajno odabranih slika za svaku od 10 klasa (3).

Klase u skupu podataka su sljedeće: avion, automobil, ptica, mačka, pas, žaba, konj, brod, kamion.



Slika 5.1: 10 slučajno odabranih slika za svaku klasu iz CIFAR-10 (3)

5.3. Rezultati

Rezultati prikazani u radu (2) predstavljaju uprosječene vrijednosti metrika. Eksperiment je ponovljen 50 puta. Metrike kojima se evaluira zaglađeni model su: standardna točnost, zaglađena točnost i certificirana točnost.

Certificirana točnost predstavlja udio točno klasificiranih primjera nad kojima djeluju neprijateljski primjeri. Koristi se za mjerenje otpornosti modela na djelovanje neprijateljskih primjera. Zaglađena točnost predstavlja udio točno klasificiranih primjera nad kojima ne djeluju neprijateljski primjeri. Koristi se za mjerenje izvedbe modela u "normalnim" uvjetima (slike bez ikakvih neprijateljskih primjera).

Standardna točnost	Zaglađena točnost	Certificirana točnost
85.53	85.15	58.5

Rezultati prikazani u radu (2)

Nakon reprodukcije ostvario sam sljedeće rezultate. Metrike nisu uprosječene pa rezultati predstavljaju ishod jednog eksperimenta.

Standardna točnost	Zaglađena točnost	Certificirana točnost
85.36	85.43	58.2

Rezultati postignuti reprodukcijom

6. Klasifikacija vrsta riba

U ovom poglavlju ću opisati provedene eksperimente sa modelom za vid zasnovanim na pažnji u problemu klasifikacije vrsta riba. Definirat ću skup podataka i njegovu pripremu, opisati parametre učenja i prikazati ostvarene rezultate.

6.1. Skup podataka

Korišteni skup podataka pripadu Institutu za oceanografiju i ribarstvo te je podijeljen na skup za učenje i skup za testiranje. Skup za učenje sadrži 1332 slike, dok skup za testiranje sadrži 717 slika. Skup podataka sadrži 3 klase: komorače iz uzgoja, divlje komorače i tune. Omjer klasa u originalnom skupu za učenje je 511 (komorače iz uzgoja) : 516 (divlje komorače) : 305 (tune). Slike su rezolucije 512 x 768. Originalni skup za učenje sam podijelio na novi skup za učenje i skup za validaciju. Skup za validaciju sadrži 33 posto slika originalnog skupa za učenje, tj. 439 slika. Sve slike su normalizirane s ImageNet parametrima. Koristio sam torchvision paket kako bih transformacijama slika iz skupa za učenje proširio skup za učenje. Transformacija nad slikama iz skupa za učenje je kompozicija sljedećih transformacija: podrhtavanje boja (color jittering), rotacija, afina transformacija. Skup za učenje u konačnici sadrži 1784 slike.

6.2. Korištene tehnologije i parametri učenja

Koristio sam Google-ov model za vid zasnovan na pažnji. Model je predtreniram na skupu podataka ImageNet. Originalni model radi sa slikama fiksne veličine 224 x 224 i veličinom isječka 16. Prilagodio sam originalni Python kod kako bi model radio sa slikama veličine 512 x 768. Također sam namjestio veličinu isječka na 64. Koristio sam Huggingface biblioteku za učitavanje, učenje i evaluaciju modela i učitavanje skupa podataka. Koristio sam torchvision paket za transformiranje slika iz skupa za

učenje. Parametri učenja su sljedeći:

- Broj epoha: 20
- Optimizator: AdamW
- Propadanje težina: 0.01
- Stopa učenja: $2 \cdot 10^{-4}$
- Veličina grupe: 8
- Dropout: 0.5
- Veličina isječka: 64

6.3. Rezultati

Metrike kojima se evaluira zaglađeni model su: točnost, preciznost, odziv i F1 mjera.

Točnost predstavlja omjer ispravno predviđenih oznaka u odnosu na ukupan broj oznaka. Mjeri koliko je ispravno predviđeno oznaka za sve primjere.

Preciznost mjeri sposobnost modela da točno predvidi pozitivne oznake za određenu klasu. Računa omjer ispravno predviđenih pozitivnih oznaka (pravi pozitivni) u odnosu na ukupan broj predviđenih pozitivnih oznaka (pravi pozitivni + lažni pozitivni).

Odziv mjeri sposobnost modela da točno identificira pozitivne oznake za određenu klasu. Računa omjer ispravno predviđenih pozitivnih oznaka (pravi pozitivni) u odnosu na ukupan broj stvarnih pozitivnih oznaka (pravi pozitivni + lažni negativni).

F1 mjera je harmonijska sredina preciznosti i odziva. Pruža uravnoteženu mjeru preciznosti i odziva, uzimajući u obzir lažne pozitivne i lažne negativne rezultate.

Točnost	Preciznost	Odziv	F1 mjera
56.21%	73.95%	56.21%	57.06%

Rezultati na skupu za testiranje (kod)

7. Zaključak

Konvolucijske i povratne neuronske mreže dominiraju područjem računalnog vida, dok model zasnovan na pažnji dominira područjem obrade prirodnog jezika. Uporaba modela zasnovanim na pažnji u računalnom vidu je i dalje ograničena. Međutim, postoje problemi u kojima je model zasnovan na pažnji primjenjiv te parira dominantim modelima u računalnom vidu. U prikazanom problemu certifikacije robusnosti modela model za vid zasnovan na pažnji se intuitivno pokazuje kao dobro rješenje za ablacije slika. Arhitektura modela za vid zasnovanim na pažnji pojednostavljuje učenje na ablacijama ignoriranjem maskiranih tokena. Također, model je brži u predikcijama od dominantnih modela računalnog vida. Očekujem da će razvoj modela za vid zasnovanim na pažnji pridonijeti boljim rezultatima certificirane robusnosti u budućnosti. Modeli zasnovani na pažnji definitivno imaju mjesto u računalnom vidu, no zasad su uglavnom ograničeni na klasifikaciju slika. Zanimljivo je promatrati njihovu primjenu u npr. detekciji objekata i segmentaciji slike.

8. Literatura

- [1] Alexander Kolesnikov Dirk Weissenborn Xiaohua Zhai Thomas Unterthiner Mostafa Dehghani Matthias Minderer Georg Heigold Sylvain Gelly Jakob Uszkoreit Neil Houlsby Alexey Dosovitskiy, Lucas Beyer. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. 2021.
- [2] Hadi Salman, Saachi Jain, Eric Wong, and Aleksander Madry. Certified patch robustness via smoothed vision transformers. 2021.
- [3] Alex Krizhevsky. CIFAR-10.
- [4] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2017.
- [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.
- [6] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016.
- [7] Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N. Dauphin. Convolutional sequence to sequence learning, 2017.
- [8] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks, 2017.
- [9] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers, 2020.

9. Sažetak

Kroz ovaj seminarski rad opisan je model zasnovan na pažnji kao dominantan model u obradi prirodnog jezika. Detaljno je opisana arhitektura modela. Opisan je i model za vid zasnovan na pažnji, njegove prilagodbe za rad sa slikama i razlike u odnosu na osnovni model zasnovan na pažnji. Nadalje, prikazana je primjena modela za vid zasnovanim na pažnji u konkretnom problemu računalnog vida. Opisan je problem ceritificirane robusnosti modela za klasifikaciju slika kao rješenje za visokorizične situacije u praksi. Opisan je zaglađeni model za vid zasnovan na pažnji koji se koristi za rješavanje spomenutog problema. Prikazani su rezultati zaglađenog model za vid zasnovnim na pažnji za vid na skupu podataka CIFAR-10 koji su predstavljeni u radu (2). Prikazana je i reprodukcija tih rezultata. Na kraju su prikazani rezultati istraživanja učinkovitosti običnog modela za vid zasnovanim na pažnji u jednom praktičnom problemu klasifikacije vrsta riba.