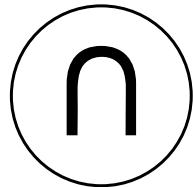


בסיסות למידה - מדריך מעשי
למהנדסים ופתחים



יסודות הלמידה העמוקה המודרנית

מהדורה 1.1 | ספטמבר 2025

ד"ר ברק אור

יסודות הלמידה העמוקה המודרנית

זה לא עוד ספר על למידה عمוקה. יסודות הלמידה העמוקה המודרנית הוא מדריך "ישומי לבנייה, אימון וIMPLEMENTATION של רשתות ניורונים عمוקות - נכתב במיוחד בברית - עבור מהנדסים ומתעניינים במערכות בעולם האמתי ולא רק בתוצאות תאורטיות.

הספר כולל Autoencoders, Transformers, LSTMs, CNNs ועוד. הוא מכסה אימון, אופטימיזציה Mixed Precision,(Lion, AdamW) ו-regularization. בסוף מוסברים Transfer Learning, Tabular Data, Time Series, Grad-CAM, SHAP Explainability .Learning התוכן תמציתי ובמיוחד הנדסי - מותאם לאנשי מקצוע. הספר משמש גם כרשומות קורס בהכשרה עצמית של ArtificialGate. כל החומר פותח כדי לתמוך בלומדים מתקעים טכנולוגיים שונים ומינים במספר שפות.

ד"ר ברק אור הוא חוקר, מרצה ויוזם בתחום הבינה המלאכותית. בעל דוקטורט בלמידה מכונה ומערכות ניוטו, ושלושה תארים מהטכניון. מייסד ArtificialGate מבית metaor.ai לעשיית הביטחוניות ולחברות המובילות בשוק הישראלי. החל מ-2024 משמש כמנהל האקדמי של המסלול לפיתוח בינה מלאכותית ולמידה عمוקה בבית הספר להייטק של Google ואוניברסיטת רייכמן.



יסודות הלמידה העמוקה המודרנית

ד"ר ברק אור

ספטמבר 2025

Artificial**Gate**

תוכן העניינים

- 4 **מבוא**
- 5 **אודות המחבר**
- 6 **1 למידת מכונה לעומת למידה عمוקה**
- 11 **2 מהי רשת נוירונים**
- 16 **3 פונקציית מחיר, Backpropagation ואופטימיזציה**
- 21 **4 איך תהליך האימון עובד?**
- 27 **5 מודי הערבת ביצועים**
- 35 **6 Regularization ו- Overfitting**
- 40 **7 מדוע אנחנו צריכים קונבולוציה?**
- 46 **8 איך פועלות רשת CNN?**
- 51 **9 רצפים וזמן: LSTM , GRU ו- RNN**
- 57 **10 הורדת ממדים עם Autoencoders**
- 63 **11 Transformer וה- Self-Attention**

68	Initialization ו Normalization 12
73	13 אוגמנטציה
78	14 אופטימיזציה מתקדמת
84	Explainability: להבין החלטות של מודלים 15
90	PyTorch מול TensorFlow 16
95	17 עבודה אפקטיבית עם Google Colab
101	Mixed Precision Training 18
106	Fine-Tuning ו Transfer Learning 19
111	20 שמירה, טעינה ו Versioning של מודלים
117	21 מימוש בסיסי
122	22 התמחויות
126	23 מפת דרכים למתנדס DL בתעשייה

מבוא

ברוכים הבאים.

ספר זה הוא תוצאה של שנים של עבודה מעשית בלמידה عمוקה - הדרכת מהנדסים, ליווי חברות והטמעת מודלים הפעילים בהיקף גדול. הוא מבוסס על הקורס היסודי שאני מלמד במסגרת ArtificialGate ונבנה בקפידה כך שייקח אתכם מאפס ועד לישום מלא.
שמי ד"ר ברק אור. אני בעל דוקטורט בלמידה מכונה למערכות ניוט. בעשור האחרון ייסדתי סטארטאפים ושיילבתי למידה عمוקה במערכות אמיתיות.

מידע נספָּך: www.barakor.com

ספר זה אינו סקירה מחקרית ואינו דיון תאורתית. הוא בסיס למהנדסים וכן עד לחת לכם:

- להבין איך למידה عمוקה עובדת מהבסיס
- לכתוב, לאמן ולהעיר מודלים ב-`PyTorch`, `Tensorflow` ב-
- לשמר ולמש מודלים בסביבות דמויות-פרודקסן
- לבחור כלים, פורמלטים ופרקטיקות עדכניות לעבודה ב-`Deep Learning`

אודות המחבר

ד"ר ברק אור הוא יzm, חוקר ומרצה המתמחה בתחום Artificial Intelligence, Machine Learning-based Inertial Sensing וב-Deep Learning.Machine Learning-based Inertial Sensing ו-Deep Learning. בעבודתו הובילה לRIESOM מספר פטנטים, לפרנסומים בכתב עת ולמחקר ישומי באקדמיה ובתעשייה. מאז 2024, ד"ר אור משמש כמנהל האקדמי של מסלול פיתוח מודלי למידה عمוקה ו-AI בבית הספר להיינט של Google ואוניברסיטת ריצ'מן. דרכו בהוראה החלה כמתרגל בטכניון וכמורה לפיזיקה בבית הספר הריאלי העברי בחיפה ובהמשך התרחבה להכשרות متقدמות בתחום DL לمهندסים ולחוקרים באקדמיה ובתעשייה.

ד"ר אור הוא מייסד ArtificialGate, פלטפורמה אינטראקטיבית מקוונת ללימוד AI, המיועדת לקהלים טכנולוגיים וארגוניים וכן יסד מספר מיזמים נוספים בתחום AI. Applied. בין אלה ניתן למנות את metaor.ai, מעבדת מחקר והכשרה; Alma Technologies, סטארט-אפ Deep-Tech עבור פתרונות מיקום לרכב מבוססי חיישנים אינרציאליים לשביבות נטולות GPS; ואת AIME Systems, העוסקת בישום AI בתחום בריאות האישה. בעבודתו בתחום אלה זכתה לمعנקים מרשות החדשנות, מפआ"ת וממשרד הכלכלה והתעשייה והניבה גם מספר פטנטים.

בשלבים מוקדמים יותר בקריירה שלו, זכה ד"ר אור בפרס גמנדר והגיע למקום השני בתחום החדשנות הביטחונית של הטכניון. הרקע האקדמי שלו כולל דוקטורט ב- Inertial Navigation-based Machine Learning מאוניברסיטת חיפה (2022), תואר M.Sc. בהנדסת אוירונואוטיקה וחיל מהטכניון (2018), תואר B.Sc. בהנדסת אוירונואוטיקה וחיל מהטכניון (2016) ותואר B.A. בכלכלת ניהול (בהצטיינות, 2016) מהטכניון.

ו שיעור 1

למידת מכונה לעומת למידה عمוקה

מבוא

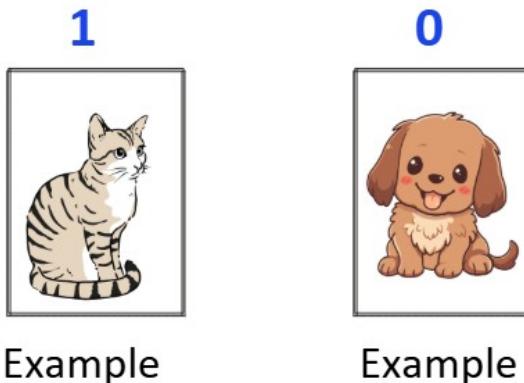
לפני שנבין איך נראה רשת נירונים, או מה פירוש Gradient Descent, חשוב לבנות בסיס: מהי בעצם למידת מכונה? מה מבדיל אותה מלמידה عمוקה? ולמה הבחנה זו קריטית למהנדסים, לחוקרים ולפתחים?

למידת מכונה: עקרונות בסיסיים

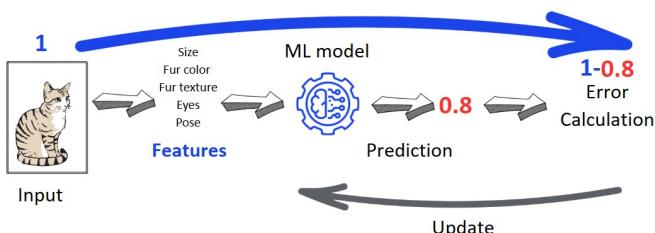
במקום לכתוב חוקים מפורשים, כמו: "אם במילilit כתוב שזcit במייל דולר - סמן אותו כ-Spam", אנו מספקים למודל אלף מיילים מתוויגים כתקינים או Spam והמודל לומד בעצמו מה מבדיל בין שתי הקטגוריות.

כל מודול זה עובר תהליכי המכונה "אימון" (Training). ב-*Supervised Learning*, אנו מספקים זוגות של קלטים ותוויות. למשל: תמונה של חתול מתוויגת כ-"1" ותמונה של כלב מתוויגת כ-"0". המודל מייצר חייזי על סמך הקלט בלבד, משווה את החיזוי לתוויות האמיתית, מחשב את השגיאה ומעדכן את עצמו כדי לסייע את השגיאה. תהליך זה חוזר על עצמו אלפי פעמים - עד שהמודל לומד את הקשר בין הקלט (התמונה) לבין הפלט הרצוי.

במהלך האימון, המודל לומד לשield תכונות (למשל צבע, טקסטורה או גובה) לתוצאה נכונה. בלמידת מכונה קלאסית עליינו להגדיר את התכונות (features) בעצמנו: קבוע מה יוזן למודל, באיזה פורמט ואיך לנормל או לאזן



איור 1.1: תמונות של חתול וכלהב, מתויגות כ-“1” וכ-“0” בהתאם.

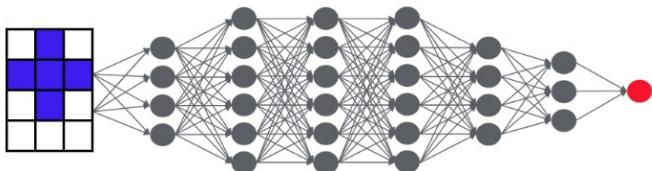


איור 1.2: תהליך האימון: חיזוי → חישוב טעות → עדכון משקלות.

את ערכי הקלט.

תפקידה של הלמידה העמוקה

כאן נכנסת לתמונה הלמידה העמוקה. למידה עמוקה אינה טכניקה ייחודית, אלא משפחה של שיטות המבוססות על רשותות נוירונים מלאכותיות. כאן המודל מקבל קלט, מעביר אותו דרך שכבות רבות של טרנספורמציה ובונה ייצוגים פנימיים - מבלי שנדרשת הגדרת פיצ'רים ידנית. אם הקלט הוא תמונה למשל, המודל מקבל את כל הפיקסלים הגולמיים ולא רק סט של פיצ'רים שהוגדרו מראש.



איור 1.3: רשת נוירונים عمוקה שמקבלת פיקסלים גולמיים כקלט.

מקרה לדוגמה: זיהוי פיצ'רים

בפרויקט יישומי מסויים, המטרה הייתה לסוג תמונות של מוצרים על בסיס פיצ'רים כמו יוורה, בהירות וניקיון. הגישה הראשונית, באמצעות למידה מכונה קלאסית, הסתמכה על מאפייני metadata מובנים - מחיר, קטגוריה ודירוגי משתמשים. התוצאה: ביצועים סבירים אך מוגבלים, כי הפיצ'רים הידניים לא היו מספיק עשיירים.

בעבר ללמידה عمוקה, אימנו מודלים ישירות על פיקסל התמונות, ללא פיצ'רים מוגדרים מראש. המודל הצליח למדוד "צוגים פנימיים" עם תבניות עדינות שלא ניתן היה להגדיר ידנית. זה הדגיש את היתרונו המרכזי של למידה عمוקה: היכולת לגלוות ולהלץ פיצ'רים רלוונטיים באופן אוטומטי, מבלי להסתמך על הנדסת פיצ'רים.

יצוג היררכי

למידה عمוקה היא תהליך של למידת "צוגים היררכית": שכבה אחר שכבה, המידע מעובד ומוצג ברמות הפשטה הולכות וגדלות - החל מפיקסלים בודדים ועד לאלמנטים מורכבים. זה מה שהופך את הלמידה העומקה לכל כך עצמאית, במיוחד עם נתונים לא-מובנים. מאז 2012, עם פריצת הדרך של AlexNet, הלמידה העומקה חוללה מהפכה בתחוםים כמו ראייה ממוחשבת, זיהוי דיבור, תרגום, ניתוח סנטימנט ויצירת תוכן.

מתי לבחור בלמידה عمוקה?

למידה عمוקה היא הבחירה הנכונה כאשר:

- מתחממים עם נתונים לא-מובנים (תמונות, טקסט, קול).
- יש כמות גדולה מאוד של נתונים זמינים לאימון.
- רוצים שהמערכת תלמד "יצוגים" בעצמה במקום להגדיר פיצ'רים באופן ידני.
- המטרה היא לגנות תכניות מורכבות שקשה או בלתי אפשרי להגדירן ידנית.

אתגרים ומוגבלות

למרות עוצמתה, ללמידה عمוקה יש אתגרים משמעותיים שיש להבאים בחשבון בישום מעשי: אימון רשותות عمוקות דורש משאבי חישוב גדולים מאוד וזמן אימון ממושך. המודלים סובלים מבעיה "הסבירתיות" (interpretability): קשה להסביר למה התקבל חיזוי מסוים, משום שהיצוגים הפנימיים מפוזרים על פני שכבות רבות ומילוני פרמטרים. מכאן נובע הינוי "black box" מערכת שבה קשה להתחקות אחרי תהליך קבלת החלטות. ובכל זאת, כאשר מאמנים ומודאים את היציבות של הרשותות, הארכיטקטורות הללו מפיקות תובנות מנוטני ענק ומגלות תכניות שהיו בלתי נגישות בשיטות מסורתיות.

סיכום

למידה عمוקה אינה רק גרסה "מתקדמת יותר" של למידה מכונה - היא שינוי תפיסתי. במקום לתקן ידנית את הפיצ'רים, אנו בונים מערכות שלומדות ליצג את הנתונים באמצעות הפשטות רב-שבכתיות.

הבדלים המרכזיים:

למידה מכונה קלאסית: השתמכות על הנדסת פיצ'רים ידנית, ארכיטקטורות פשוטות, נתונים מובנים.

למידה عمוקה: למידת פיצ'רים אוטומטית, שימוש בארכיטקטורות עמוקות,
יעילה במיוחד עבור נתונים לא-מובנים.
בשיעורים הבאים נלמד איך לבנות מודלים מסוג זה - החל מהנוירון
המלאכותי הבודד ועד לרשותות עמוקות שלמות.

[שיעור ראשון בקורס זמין ב-YouTube](#)

Unit 2 שיעור 2

מהי רשת נוירונים

מבוא

בשיעור הקודם בחנו את ההבדל בין למידת מכונה למידה عمוקה. ראיינו של מידה عمוקה מאפשרת למודל ללמידה ייצוגים פנימיים בעצמו, ללא הנדסת פיצ'רים ידנית. כתה עמוקה במבנה הפנימי של רשתות נוירונים מלאכותיות, בניסוח המתמטי שלחן ובסיבות שבזוכותן הן מסוגלות לייצג פונקציות מורכבות.

הנוירון המלאכותי

בבסיס כל רשת נוירונים נמצא נוירון מלאכותי, או Perceptron. הנוירון מקבל וקטור קלט $\mathbf{x} \in \mathbb{R}^n$, מכפיל כל רכיב במשקלות מתאימה, מסכם את התוצאות, מוסיף מזוזיף את הסכום $b \in \mathbb{R}$ bias ו מעביר את הסכום דרך פונקציית אקטיבציה φ כדי לייצר פלט :

$$(2.1) \quad y = \varphi(\mathbf{w} \cdot \mathbf{x} + b)$$

כאשר:

$\mathbf{x} = [x_1, x_2, \dots, x_n]^\top$ • הוא וקטור הקלט.

$\mathbf{w} = [w_1, w_2, \dots, w_n]^\top$ • הוא וקטור המשקלות.

- $x \cdot w$ היא המכפלה הסקלרית $\sum_{i=1}^n w_i x_i$.
- b הוא ערך ה-bias.
- φ היא פונקציית אקטיבציה לא-יליניארית.
- y הוא הפלט הסקלרי של הנeuron.

למרות שנeuron בודד הוא פונקציה פשוטה, צירוף של רבים מהם לשכבות מאפשר למודל לייצג פונקציות מורכבות מאוד.

פונקציות אקטיבציה

פונקציית האקטיבציה φ מוסיפה אי-יליניאריות למודל. ללא אי-יליניאריות גם אם נערום שכבות רבות, נקבל בסופו של דבר טרנספורמציה ליניארית אחת - שאינה מספקת לביעות מורכבות. שתי פונקציות נפוצות:

:ReLU (Rectified Linear Unit) •

$$(2.2) \quad \varphi(x) = \max(0, x)$$

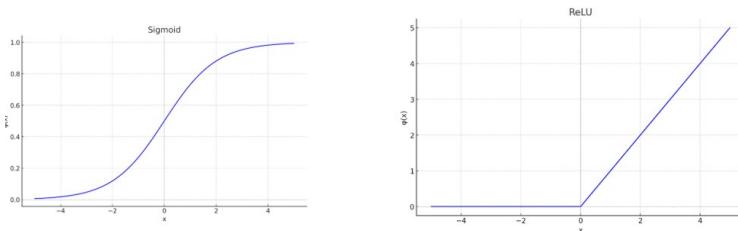
מחזירה 0 עבור ערכים שליליים ומשAIRה את הערכים החיוביים כפי שהם. פונקציה זו פשוטה, יעה חישובית ונפוצה מאוד ברשומות מודרניות.

:Sigmoid •

$$(2.3) \quad \varphi(x) = \frac{1}{1 + e^{-x}}$$

דוחשת את הערכים לטוווח (0, 1). שימושה רבות ברשומות מוקדמות, במיוחד במשימות סיוג ביןארי, אך כיום השימוש בה מועט בשל בעיות .Vanishing Gradient

טייפ: אם פונקציית המחוור לא יורצת כמהלך האימונו, אחת הסיבות האפשרות היא בחירה לא מתאימה של פונקציית אקטיבציה - למשל שימוש כ- *Sigmoid* *ReLU* בשכבות חכויות במקומות



איור 2.1: פונקציות אקטיבציה: Sigmoid (משמאל) וReLU (מימין).

מבנה שכבות ברשת

רשת נוירונים عمוקה מורכבת ממספר רב של שכבות. כל שכבה מקבלת את הפלט של השכבה הקודמת ומחשבת פלט חדש:

$$(2.4) \quad \mathbf{y}^{(\ell)} = \varphi(\mathbf{W}^{(\ell)}\mathbf{y}^{(\ell-1)} + \mathbf{b}^{(\ell)})$$

כאשר:

- $\mathbf{y}^{(\ell-1)}$ הוא הפלט לשכבה ℓ (כלומר פלט השכבה הקודמת).
- $\mathbf{W}^{(\ell)}$ היא מטריצת המשקלות של שכבה ℓ .
- $\mathbf{b}^{(\ell)}$ הוא וקטור ה- bias של שכבה ℓ .
- φ מופעלת על כל רכיב בנפרד.

הארQUITטורה הבסיסית כוללת:

- שכבת קלט - מקבלת את הנתונים הגולמיים.
- שכבות חבויות - לומדות ייצוגים פנימיים מורכבים.
- שכבת פלט - מייצרת תחזית בהתאם למשימה.

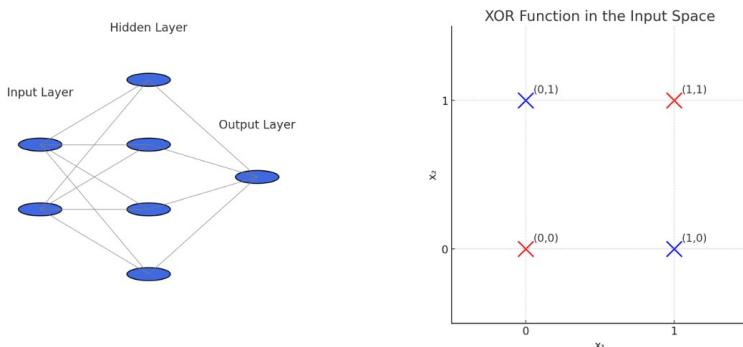
רַב-שְׁכָבִי (MLP) Perceptron

(Multi-Layer Perceptron) MLP הוא סוג של רשת שבה כל שכבה מחוברת במלואה לשכבה שאחריה.

יתרונו מרכזי: MLP מסוגל לייצג בעיות שאינן ניתנות להפרדה ליניארית.
דוגמה קלאסית היא פונקציית XOR:

$$(2.5) \quad y = \begin{cases} 1 & x_1 \neq x_2 \\ 0 & \text{אחרת} \end{cases} \quad x_1, x_2 \in \{0, 1\}$$

פונקציה זו אינה ניתנת להפרדה על ידי קו ישר במרחב הקלט. אך MLP עם שכבה חבויה קטנה אחת בלבד מסוגל ללמידה אותה בהצלחה.



איור 2.2: סיווג XOR: משמאל - רשת נוירונים; מימין - גרף של פונקציית XOR.

משפט הקירוב האוניברסלי

The Universal Approximation Theorem states that a neural network with a single hidden layer can approximate any continuous function on a compact domain. This means that if you have a function $f: D \rightarrow \mathbb{R}$ where D is a compact subset of \mathbb{R}^n , then there exists a neural network \hat{f} such that $\hat{f}(x) \approx f(x)$ for all $x \in D$. The number of neurons in the hidden layer depends on the complexity of the function f and the desired accuracy of the approximation. The proof of this theorem relies on the fact that the hidden layer can represent any continuous function on a compact domain, and the output layer can map the hidden layer's representation to the final output. This allows the network to learn complex patterns and relationships in the data.

סיכום

רשת נוירונים היא פונקציה פרמטרית הבנויה משכבות של נוירונים מלאכותיים. כל נוירון מחשב צירוף ליניארי משוקל של הקלט, מוסיף bias ו מעביר את התוצאה דרך פונקציית אקטיבציה לא-لينיארית. באמצעות צירוף של שכבות רבות, הרשת מסוגלת ללמוד ייצוגים מורכבים ולהתמודד עם בעיות שלא ניתנות לפתרון באמצעות טרנספורמציה ליניארית בלבד. זהו הצעד הראשון לעבר רשותות עמוקות שambilאות לידי ביטוי את מלא העוצמה של הלמידה העמוקה.

שיעור 3

פונקציית מחיר, Backpropagation ואופטימיזציה

מבוא

בשיעור הקודם רأינו כיצד מידע עבור קדימה מהקלט לפט ברשת נירונים כדי ליחסו. אבל רשות אינה "נולדת" עם ידע - עליה ללמידה מהנתונים. בשיעור זה נבין כיצד מתרחש תהליכי הלמידה, באמצעות שלושה מרכיבים מרכזיים: פונקציית המחיר, Backpropagation ושיטת האופטימיזציה Gradient Descent.

פונקציית מחיר

פונקציית המחיר (Loss Function) מודדת עד כמה החיזוי של המודל רחוק מהתאוג האמתי. המטרה: לכמת את הפער בין הפלט החזוי לבין התשובה הנכונה.

כי נניח $y \in \{0, 1\}$ היא התווית האמיתית (התאוג) ו- $\hat{y} \in (0, 1)$ היא ההסתברות שהמודל חזה. פונקציה נפוצה מאוד לשיווג ביןاري היא Cross-Entropy:

$$(3.1) \quad L = - \left[y \cdot \log(\hat{y}) + (1 - y) \cdot \log(1 - \hat{y}) \right]$$

אם החיזוי קרוב לערך הנכון $\leftarrow L$ קטן. אם החיזוי רחוק או שגוי $\leftarrow L$

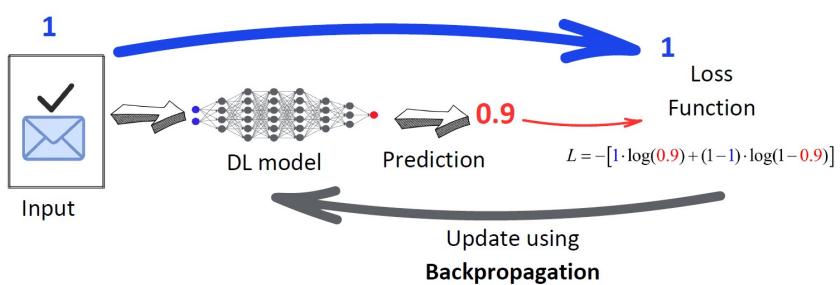


איור 3.1: דוגמה: מייל אחד מסומן כ-Spam והשני לא. המשימה - לחזות נכון. נכון.

גודל.

תרשים זרימה מלא של אימון

תהליך האימון מתבצע באופן הבא:
קלט ← Backpropagation ← חישוב טעות ← חיזוי ← עדכון משקولات.



איור 3.2: תרשים זרימת האימון.

Backpropagation

כדי לשפר את המודל, علينا לעדכן את המשקלות כך שיקטינו את הערך של פונקציית המחיר. Backpropagation הוא האלגוריתם שמאפשר זאת באמצעות חישוב נגזרות ושימוש בכלל השרשרת. הגרדיאנט של פונקציית המחיר ביחס לכל משקל w מוגדר כך:

$$(3.2) \quad \frac{\partial L}{\partial w}$$

באמצעות כלל השרשרת (Chain Rule), מחושבים הגרדיאנטים מהשכבה האחרוןה אחורה - דרך שכבות הביניים - עד לקלט. כך ניתן לשיך לכל פרמטר (משקלות) במודל את התרומה שלו לשגיאה ולעדכן אותו בהתאם.

הגרדיאנט - אינטואיציה

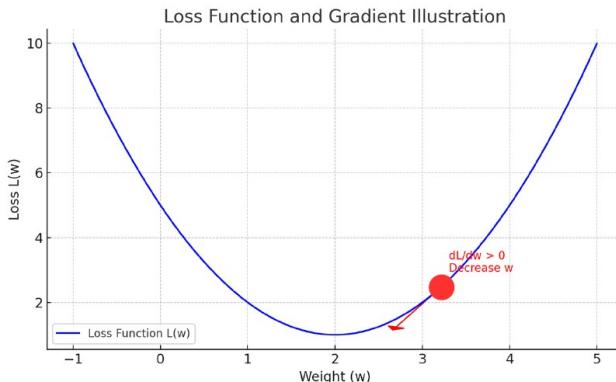
נניח שפונקציית המחיר $L(w)$ היא פונקציה חלקה של משקל יחיד w . הנגזרת בנקודת מסויימת מתארת את השיפוע באותה נקודה: ככלומר, נוכל לקבוע אם הגדלת w תעלה או תקטין את L .
עדכן המשקלות מתבצע לפי הכלל:

$$(3.3) \quad w_{k+1} = w_k - \eta \cdot \frac{dL}{dw}$$

כאשר $0 < \eta$ הוא קצב הלמידה (Learning Rate).

קצב הלמידה

קצב הלמידה η קובע את גודל הצעד בכל עדכו:
 - אם η קטן מדי - המודל יתכנס לאט מאוד.
 - אם η גדול מדי - המודל עלול להתבדר, ככלומר הטעות תגדל במקומות רחוקו.



איור 3.3: פונקציית המחיר כתלות במשקל. החץ האדום מראה את כיוון הגרדיינט. העדכון נע בכיוון של הקטנת טעות החיזוי - מזעור פונקציית המחיר.

טיוף: אם במהלך האימון נראה שהטעות " קופצת" למעלה ולמטה ולא מתכנסת לערך יציב, סיבה אפשרית היא שקצב הלמידה גבוהה מדי.

כפי שתיאר זאת אחד הסטודנטים: "יוטיוטי��צ' למיזה של 1 והמזהל פשוט קפץ הלוֹץ ושוב כמו כדור פינג-פונג."

דוגמה: מודל ליניארי מופשט עם משקלות יחידה

נבחן מודל עם משקלות יחידה w וקלט x :

$$(3.4) \quad \hat{y} = w \cdot x$$

פונקציית המחיר - טעות ריבועית ממוצעת (MSE) עבור דוגמה אחרת:

$$(3.5) \quad L = (\hat{y} - y)^2$$

נניח: $x = 2, y = 6, w = 1$

$$\hat{y} = 1 \cdot 2 = 2$$

הטעות: $L = (2 - 6)^2 = 16$

נחשב את הגרדיינט:

$$(3.6) \quad \frac{dL}{dw} = 2 \cdot (wx - y) \cdot x$$

נציב $w = 1, x = 2, y = 6$

$$\frac{dL}{dw} = 2 \cdot (2 - 6) \cdot 2 = -16$$

עדכון עם $\eta = 0.1$:

$$(3.7) \quad w_{\text{new}} = w - \eta \cdot \frac{dL}{dw} = 1 - 0.1 \cdot (-16) = 2.6$$

כעת, לאחר עדכון ייחד, המשקל w השתפר משמעותית והתקרב לערך שיקטין את השגיאה.

סיכום

תהליך הלמידה של רשת נוירונים מורכב משלושה שלבים מרכזיים:

1. חישוב החיזוי (Forward Pass).
2. חישוב פונקציית המחיר - מדידת המרחק מהתוצאה הנכונה.
3. חישוב גראדיינטים באמצעות Backpropagation ועדכון משקלות בעזרת Gradient Descent.

שלשות המרכיבים הללו הם לבת תהליכי הלמידה ברשומות עמוקות.

4 שיעור

איך תהליכי האימון עובדים?

עד כה בחנו כיצד רשתות נוירונים בנויות וכיצד הן מייצרות חיזויים. אך הLIBHE של כל אלגוריתם למידה טמונה ביכולת שלו להשתפר. איך בדיקת רשת לומדת מתוך DATA? מה מתרחש במהלך האימון וכיצד הפרמטרים - המשקولات וההתווית - מתעדכנים כדי לצמצם את השגיאה?

לעומת Mini-Batch Gradient Descent Batch

בצורתו הקלאסית, הידועה כ-**Batch Gradient Descent**, המודל מעבד את כל DATAהסט האימון בעבר קדימה אחד, מחשב את פונקציית המחיר עבור כל הדוגמאות ומבצע ממוצע כדי לקבל ערך סקלרי יחיד:

$$(4.1) \quad L_{\text{batch}} = \frac{1}{N} \sum_{i=1}^N \ell(y_i, \hat{y}_i)$$

כאן N הוא מספר הדוגמאות הכוללות באימון, y_i הוא הליביל האמתי, \hat{y}_i הוא החיזוי של המודל עבור הקלט x_i ו- ℓ היא פונקציית המחיר לכל דוגמה. הגרדיינט של פונקציית המחיר ביחס לפרמטרי המודל w מחושב:

$$(4.2) \quad \nabla_w L_{\text{batch}} = \frac{1}{N} \sum_{i=1}^N \nabla_w \ell(y_i, \hat{y}_i)$$

לבסוף, המשקولات מתעדכנות באמצעות צעד גראדיינט יחיד:

$$(4.3) \quad w_{k+1} = w_k - \eta \cdot \nabla_w L_{\text{batch}}$$

שיטת זו מספקת הערכה לגרדיינט, יציבה ומודוקת, אך אינה יעילה חישובית ודורשת זיכרון רב עבור נתונים גדולים. לכן, היא כמעט ואינה בשימוש בפועל.

במקום זאת, משתמשים בדרך כלל ב-**Mini-Batch Gradient Descent**, שבו הדאטסהסט מוחלק לקבוצות קטנות יותר בגודל B , הנקראות **Mini-Batches**. עבור כל מיני-באץ', המודל מבצע את השלבים הבאים:

1. מעבר קדימה על כל B הדוגמאות לחישוב החיזויים $\hat{y}_1, \dots, \hat{y}_B$

2. חישוב פונקציית המחיר עבור כל דוגמה $\ell(y_j, \hat{y}_j)$

3. ממוצע החישוביים כדי לקבל את פונקציית המחיר של המיני-באץ':

$$(4.4) \quad L_{\text{mini}} = \frac{1}{B} \sum_{j=1}^B \ell(y_j, \hat{y}_j)$$

4. חישוב הגרדיינט של פונקציית המחיר:

$$(4.5) \quad \nabla_w L_{\text{mini}} = \frac{1}{B} \sum_{j=1}^B \nabla_w \ell(y_j, \hat{y}_j)$$

5. עדכון המשקלות:

$$(4.6) \quad w_{k+1} = w_k - \eta \cdot \nabla_w L_{\text{mini}}$$

חשוב: המשקלות מתעדכנות פעמי אחת לכל מיני-באץ', רק לאחר חישוב הממוצע והגרדיינט. גישה זו מאפשרת שימוש יעיל יותר במשאבים חישוביים ומספקת תהליכי אופטימיזציה יציב יותר בהשוואה לעדכון על כל דוגמה בודדת.

Updates , Iterations , Epochs

כדי להבין כיצד האימון מתקדם לאורך זמן, מגדירים מספר מונחים מקובלים:

- **Epoch** - מעבר מלא אחד על כל DATASET האימון.
- **Mini-Batch** - תת-קובוצה מתוך DATASET האימון (בגודל B) המשמשת לעדכון משקלות אחד.
- **Iteration** - צעד עדכון משקלות אחד, לאחר עיבוד מיני-באץ' יחיד.

מספר האיטרציות בכל Epoch נקבע על פי:

$$(4.7) \quad \text{Iterations-per-epoch} = \left\lceil \frac{N}{B} \right\rceil$$

דוגמא: עבור $N = 10,000$ דוגמאות אימון וגודל מיני-באץ' 32
מספר האיטרציות בכל Epoch הוא:

$$(4.8) \quad \left\lceil \frac{10,000}{32} \right\rceil = 313$$

משמעות הדבר שהמודל מבצע 313 מעברים קדימה ואחוריה בכל Epoch, כל אחד מלוחה בעדכון משקלות. רשתות נירוניות מאומנות לרוב על פניהם ששרותות ואף מאות Epochs כדי להגיע להתקנסות.

אם פונקציית המחיר אינה יורצת לאחר מספר Epochs, ייתכו שהלמידה הגיעה לרווחה ושיעזרו או ייכילו מחדש את הפרמטרים.

Learning Rate

כפי שראינו קודם, בעוד שהגרדיינט מספק את הכיוון שבו על המודל לשנות את משקלותיו, ה-**Learning Rate** η קובע את גודל הצעד בכיוון זה. אם η קטן מדי - ההתקנסות איטית. אם הוא גדול מדי - תהליך האימון עלול להיות לא יציב, להתנווד או אף להתבדר.

הגרדיינט מצביע על הזרץ - אך ה-**Learning Rate** η מוכיח כמה ניל.

בחירת Learning Rate מתאימים היא קריטית לאימון יעיל ויציב.

Optimizers

ה-Optimizer הוא האלגוריתם שאחראי על יישום עדכוני הגרדיינט לפרמטרי המודל. הוא משתמש בגרדיינטים שחושבו באמצעות Backpropagation, מכפיל אותם ב- η ומבצע עדכון משקלות. הזרה פשוטה ביותר היא :Gradient Descent

$$(4.9) \quad w_{k+1} = w_k - \eta \cdot \nabla_w L_k$$

בפועל, שיטות עובדה מודרניות לאימון עושות שימוש באופטימיזרים מתקדמים יותר כגון Adagrad, Adam או RMSProp, אשר משלבים מנגנוןים לייצוב והאצת ההתקנסות.

שיעור 14 יעמיק במבנה המתמטי וcite{cite} וחרוגות של אופטימיזרים אלו.

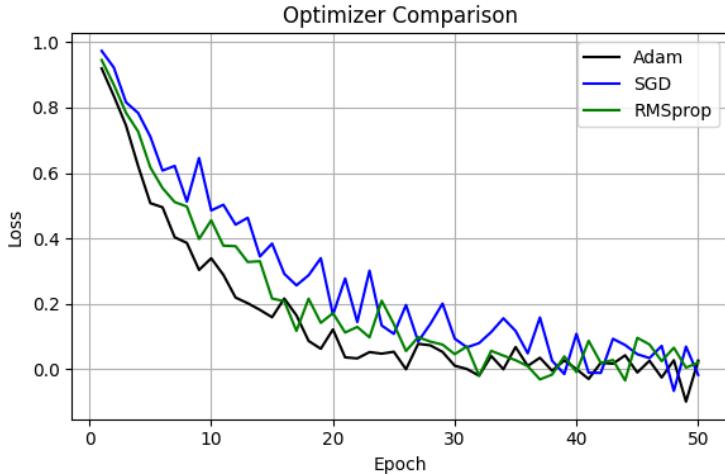
המחשה של דינמיקת האימון

כדי לבחון כיצד היפר-פרמטרים המשפיעים על ההתקנסות, ניתן לדמות אימון מודל על דאטסהט סינטטי תוך שינוי האופטימיזר, ה-Learning Rate וגודל ה-Batch.

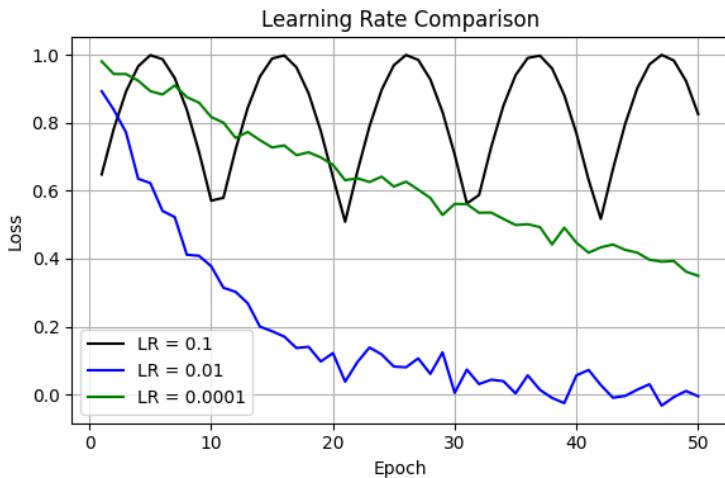
התנהגות Optimizer: Adam (שחור) מתקנס במהירות ובצורה חלקה. SGD (כחול) מציג חוסר יציבות. RMSProp (ירוק) יציב יותר מ-Adam פחות יעיל מ-Adam.

התנהגות Learning Rate: קצב למידה גבוה (שחור) מוביל לעדכוניות רועשים או מתבדרים. קצב נמוך (ירוק) מאט את ההתקנסות. ערך בינוני (כחול) מספק איזון בין יציבות ללמידה.

התנהגות Batch Size: באצ'ים קטנים (שחור) מייצרים עדכונים מהירים אך רועשים. באצ'ים בינוניים (כחול) מספקים איזון. באצ'ים גדולים (ירוק) מובילים להתקנסות איטית אך חלקה יותר.



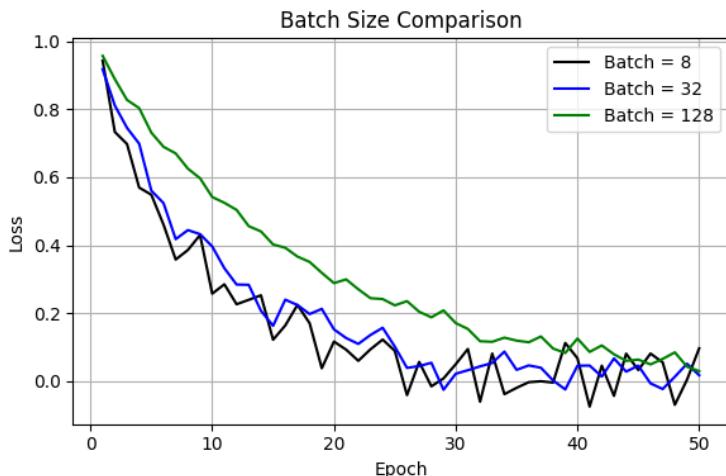
אייר 4.1: השוואת השוואת Loss - Optimizers עלות Epoch. שחור: Adam, כחול: RMSProp, ירוק: SGD



אייר 4.2: השוואת השוואת Learning Rate. גובה: גבוהה, כחול: בינוני, ירוק: נמוך.

סיכום

динמיקת האימון נקבעת על ידי האינטראקציה בין מספר רכיבים מרכזיים.
מודל ה-**Batch** מכתיב כמה דוגמאות מעובדות יחד לצורך חישוב עדכון $\frac{1}{n}$



איור 4.3: השוואת Batch Size. שחור: קטן, כחול: בינוני, ירוק: גדול.

גרדיינט ומשפיע על היציבות והווריאנס של תהליכי האופטימיזציה. **מספר Epochs** קובע כמה פעמים המודל עובר על כל>Dataהסט האימון ומשפיע על החשיפה הכלולת לDATA. **Learning Rate** מגדיר את גודל עדכון במרחב הפרמטרים ומażן בין מהירות התוכניות לבין סיכון של פספוס. לבסוף, **Optimizer** מגדיר את חוק העדכון בפועל ולעיתים משלב מנגנים כמו Momentum או התאמות אדרטיביות לשיפור יעילות הלמידה.

אימון אפקטיבי תלוי בבחירה נכונה של נתונים לכל אחד מהרכיבים הללו. גם טעויות קטנות עלולות לגרום לעדכונים לא יציבים, להתקנסות איטית, או לכישלון בלמידה.

בשיעור הבא נעבור מאופטימיזציה להערכת ונבחן כיצד לקבוע אם המודל אכן למד ומסוגל להקליל על DATA שלא ראה קודם.

שיעור 5

מדדי הערכת ביצועים

מבוא

בשיעורים הקודמים בנו רשת נוירונים, אימנו אותה על דатаה אמיתי, עדכנו משקלות באמצעות גרדיאנטים ובחרנו אופטימיזר. בעת מגיעה השאלה הクリיטית: איך נדע אם המודל שלנו באמת טוב? המטריה של Machine Learning אינה לשנן את דатаה האימון, אלא להכליל - ככלומר, להצליח ביצועים טובים גם על דатаה חדש שלא נראה קודם. יכולת זו נקראת **יכולת ההכללה** של המודל והיא המדריך המרכזי של כל מערכת חישוי.

פיקול דעתהסט

לפני שמתחליל כל אימון - אפילו לפני בחירת הארכיטקטורה של המודל - הדעתהסט חייב להיות מחולק לחת-קבוצות נפרדות. חלוקה זו חיונית כדי לוודא שביצועי המודל משקפים במידה אמיתי ולא שינון. ה-**Training Set** משמש להתאמת המודל: חישוב הגרדיאנטים מtabצע עלייו והפרמטרים מתעדכנים בהתאם. ה-**Validation Set** נשמר לצד במהלך האימון ומשמש לניטור ביצועים, כיוון היפר-פרמטרים זיהוי סימני Overfitting. לבסוף, ה-**Test Set** נשמר אך ורק להערכת ההכללה הסופית - אין להשתמש בו בשום שלב של בחירת מודל או כיול. אסטרטגיה נפוצה מחלוקת בערך 70% מהדатаה לאימון, 15% לוילדציה

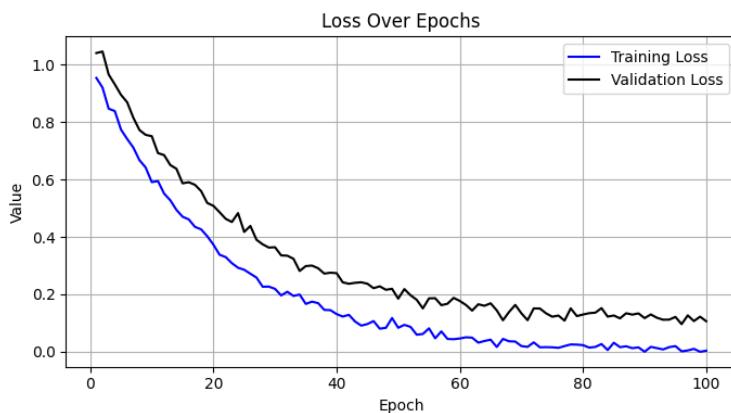
�-15% לבדיקה. עם זאת, היחסים עשויים להשנות בהתאם לגודל ולאופי הדאטסהט. העיקנון המרכזי הוא שה-Test Set חייב להישאר בלתי נגש עד הסוף, כדי לספק הערכה אובייקטיבית לביצועים אמת.

מעקב אחרי Accuracy ו-Loss

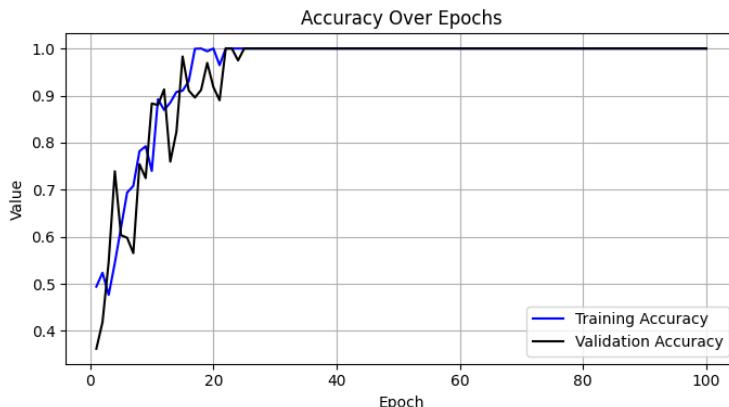
במשימות סיוג עוקבים אחרי שני מדדים מרכזיים לאורך האימון: **Loss** ו-**Accuracy**. מחושבים בנפרד על דאטסהט האימון ועל דאטסהט הולידייטה. פונקציית ה-**Loss** מספקת אותן רציף וגזר המראה עד כמה רוחקות תוצאות המודל מהלייבלים האמיתיים. היא רגישה לביטחון התוצאות ומענישה יותר על שגיאות גדולות. לעומת זאת, **Accuracy** הוא מדד בודד, הבודק את שיעור התוצאות הנכונות:

$$(5.1) \quad \text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}}$$

שימוש לב: ערך נמוך של Loss לא בהכרח מעיד על Accuracy גבוה - ה-**Loss** בוחן את ביטחון התוצאות, בעוד **Accuracy** רק לנוכנות התוצאות.



איור 5.1: מחיר (Loss) באימון ובולידייטה לאורך Epochs: כחול - אימון, שחור -olidיטה.



איור 5.2: דיק (Accuracy) באימון ובוילידציה לאורך Epochs: כחול - אימון, שחור -olidציה.

Confusion Matrix ודוגמה

כדי להבין טוב יותר את התנהוגות המודל בסיווג בינארי, במיוחד בתחוםים רגשיים כמו בריאות, משתמשים בכלי בשם **Confusion Matrix**. הוא משווה תוצאות עם ליבלים אמיתיים ומחלק את התוצאות לארבעה מצבים:

Confusion Matrix

Actual Negative	Actual Positive	
False Positive (FP)	True Positive (TP)	Predicted Positive
True Negative (TN)	False Negative (FN)	Predicted Negative

הבדיקה (Accuracy) במונחי מרכיבי הטבלה:

$$(5.2) \quad \text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

דוגמה: זהויות מחלה

נניח סיוג ביןארי לגילוי מחלת נדירה שבה רק אחד מכל 100 אנשים חולה. אם המודל ינבא "בריא" לכולם, הוא ישיג 99% Accuracy - מספרית גבוהה, אך מעשית קטסטרופתית. מטריצת הבלבול חשפת מקרים מטעים כאלה.

Healthy Predicted:	Sick Predicted:	
(Missed) FN	(Correct) TP	Sick Actually
(Correct) TN	Alarm) (False FP	Healthy Actually

לכל רביעון במטריצה יש השלכות בעולם האמיתי:
TP המחלה זוהתה נכון.

FN חולה הוחמצ - כשל קרייטי.
FP בריא זוהה כחולה - יוצר חרדה ועלות מיותרת.
TN בריא סוווג נכון.

בהתדרות כאלה, ה-"**חיובי**" מתייחס למצב שאותו רוצים לגלות - כאן, חוליו. כל המדדים הבאים (**Recall**, **Precision**) מחושבים ביחס למחלוקת החיויבות.

F1 Score ו- Recall ,Precision

מהו כך המטריצה מדדים אינפורטטיביים יותר:
Precision - איזה חלק מהתחזיות החיויבות נכון?

$$(5.3) \quad \text{Precision} = \frac{TP}{TP + FP}$$

- איזה חלק מהמקרים החיויבים זוהה? **Recall (Sensitivity)**

$$(5.4) \quad \text{Recall} = \frac{TP}{TP + FN}$$

:**Recall** ו- **Precision** של **F1 Score** - המוצע ההרמוני של

$$(5.5) \quad F_1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

מדדים אלה חשובים במיוחד בדאטאסתים לא מאוזנים. Precision גבוה עם Recall מuido על מודל שמרני; Precision גבוה עם Recall נמוך מרמז על הרבה התראות שווה.

AUC ו- ROC Curve

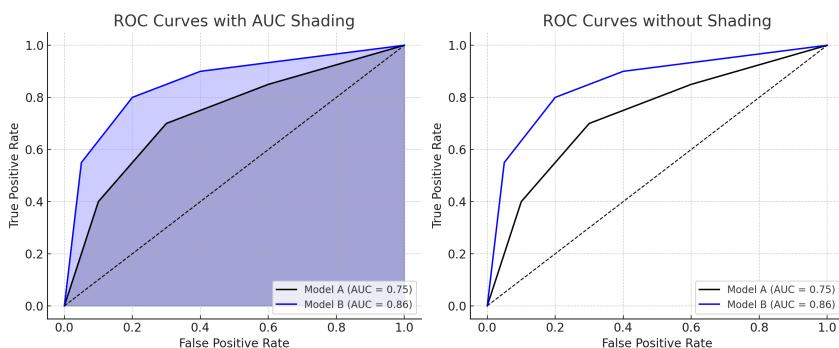
כאשר מודלים מחזירים הסתבריות, מגדרים סף החלטה (בדרכן כלל 0.5) כדי להמיר לSieog ביןארי. שינוי הסף משפייע על רגישות וספכיפיות.

עקומת ROC מטארת את True Positive Rate (Recall) מול False Positive Rate עבור ספים שונים. השיטה שמתוחת לעוקומה - ה-AUC - מסכם את יכולת ההבחנה של המודול בין מחלקות:

• Sieog מושלם : 1.0 = AUC •

• לא טוב מניחוש אקראי : 0.5 = AUC •

• תחזיות היפות : $0.5 < \text{AUC}$ •



איור 5.3: עקומת ROC עם AUC (שטח מוצלל): מודל A - שחור, מודל B - כחול, $\text{AUC}=0.75$; $\text{AUC}=0.86$.

מדדי רגסיה

כאשר המודל מוחזיר ערך רציף ולא מחלקה, משתמשים במדד **Regression**alo מודדים את הפער בין ערכים חזויים \hat{y}_i לערכים אמיתיים y_i .

Mean Absolute Error (MAE)

$$(5.6) \quad \text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

MAE מודד את גודל השגיאה הממוצע, מתייחס באותו אופן לכל סטייה. הוא חסין למדידות חריגות ובעל ייחדות זהות למשתנה החזוי.

Mean Squared Error (MSE)

$$(5.7) \quad \text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

MSE מדגיש שגיאות גדולות בגלל העלאה בריבוע. נפוץ כמשמעותם בו כפונקציית מחיר כי הוא רציף ונזר, אך רגש יותר לחריגים.

מקדם ההסביר (R^2)

ה- **R^2** , **Coefficient of Determination**, המסומן R^2 , הוא מדד נפוץ להערכתה מודלי רגסיה. הוא מודד עד כמה החיזויים \hat{y}_i קרובים לערכים האמיתיים y_i בהשוואה לביסולין פשוט - ניבוי תמידי של הממוצע.

פורמלית:

$$(5.8) \quad R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

כאשר SS_{res} הוא סכום ריבועי השגיאות בין חיזויים לערכים אמיתיים ו- SS_{tot} הוא סכום הוריאציה בערכים האמיתיים סביב הממוצע \bar{y} :

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

$$SS_{\text{res}} = \sum_{i=1}^n (y_i - \hat{y}_i)^2, \quad SS_{\text{tot}} = \sum_{i=1}^n (y_i - \bar{y})^2$$

ערך $R^2 = 1$ מצביע על התאמה מושלמת. ערך $R^2 = 0$ מצביע על חיזוי גרוע כמו הממוצע. אם $0 < R^2 < 1$, המודל גרוע מהביזטליון.

מדדים נוספים לסיווג

Balanced Accuracy

מדד זה מתחשב בדאטסהטים לא מאוזנים על ידי ממוצע ה-Recall בין מחלקות:

$$(5.9) \quad \text{Accuracy Balanced} = \frac{1}{2} \left(\frac{TP}{TP + FN} + \frac{TN}{TN + FP} \right)$$

כך כל מחלוקת משפיעה באופן שווה, בלי קשר לשכיחות.

Matthews Correlation Coefficient (MCC)

$$(5.10) \quad \text{MCC} = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

ה-MCC הוא מדד מאוזן גם מול דאטה לא מאוזן. מתחשב בכל ארבע התוצאות ומתרנגן כמדד מתאם: $+1 =$ תחזית מושלמת, $0 =$ ניחוש אקראי, $-1 =$ חוסר התאמה מוחלט.

סיכום

הערכת ביצועי מודל דורשת יותר ממספר יחיד. בסיווג, ה- Accuracy הוא נקודת פтиחה, אך עמוק יותר מגע מ- Precision, Recall, F1 Score ו- AUC. ברגression, MSE, MAE ו- R^2 מודדים את איכות התוצאות המספריות. מדדים כמו MCC Balanced Accuracy מספקים רובסטיות בסביבות לא מאוזנות. מעל הכל, הערכה חייבת להתבצע על Test Set שמור כדי לספק הערכת הכללה אמינה.

שיעור 6

Regularization ו- Overfitting

מבוא

בשלבים המוקדמים של אימון רשת נירוניים, הכל עשויה להיראות תקין: ה- Training Loss יורד בהדרגה וה- Accuracy משתפר מאפקט לאפקט. כאשר מוצג דатаה חדש ולא נראה, ביצועי המודל עלולים לקרוס בפתאומיות. תופעה זו אינה באג - זהה בעיה יסודית בלמידה عمוקה, הידועה בשם Overfitting (התאמת יתר).

Overfitting מתרחש כאשר המודל מתאים עצמו יותר מדי לדאטatest האימון, קולט לא רק את הדפוסים הבסיסיים אלא גם רעשים, חריגים וקורלציות מקרים. כתוצאה לכך, המודל מצליח על הדטה שראתה, אך נכשל בהכללה על דוגמאות חדשות. לעומת זאת, Underfitting (תת התאמת), מתרחש כאשר המודל אינו מצליח לקלוט את הדפוסים בדאטatest האימון, לרוב בשל קיבולת לא מספקת או אימון לא מספק. גם Underfitting ו- Overfitting פוגעים ביכולת ההכללה - שהיא המטרה המרכזית.

דוגמה

דמיינו תלמיד ש shinon תשובה מבחן בשנה שעברה מבלי להבין את החומר. כאשר מוצג לו מבחן חדש עם שאלות מעט שונות, הוא אינו מסוגל לענות

נכון. האנלוגיה פשוטה: שינון איננו למידה. באופן דומה, רשת נוירונים עשויה להמשיך להמשיך להקטין את Training Loss בעוד Validation Loss מתחילה לעלות - סימן קלאסי ל-Overfitting. בתחילת המודל משתפר גם על אימון וגם על ולידציה, אך בשלב מסוים הוא מתחילה לשינוי במקומות ללימוד.

התנוגות ה-Loss ככלי אבחון

במהלך האימון עוקבים אחרי ה-Loss גם על דاطה הסט האימון וגם על הולידציה. נסמן $L_{\text{val}}(t)$ ו- $L_{\text{train}}(t)$ כמחיר באימון ובולידציה, בהתאם, באפקט t .
בשלבים המוקדמים:

$$L_{\text{train}}(t) \downarrow, \quad L_{\text{val}}(t) \downarrow$$

:Overfitting בנקודת תחילת

$$L_{\text{train}}(t) \downarrow, \quad L_{\text{val}}(t) \uparrow$$

הסתירה הזאת היא סימן קריטי לכך שהמודל כבר לא לומד מבנה משמעותית, אלא משנה את דاطה הסט האימון.

מדוע מתרחש Overfitting

רשתות נוירוניים מודרניות מכילות לעיתים מיליוני פרמטרים. עם מספיק קיבולת, המודל יכול תמיד להשיג Training Loss כמעט אפסי - אפילו על דاطה עם תוויות אקריאיות. אך פתרון זה לא לומד רושם במקומות מבנה בסיסי. דיקוק מושלם על דاطה הסט האימון הוא לרוב דגל אדום, לא סיבה לחגינה. הוא מצביע על עודף פרמטרים ועל היעדר Inductive Bias.

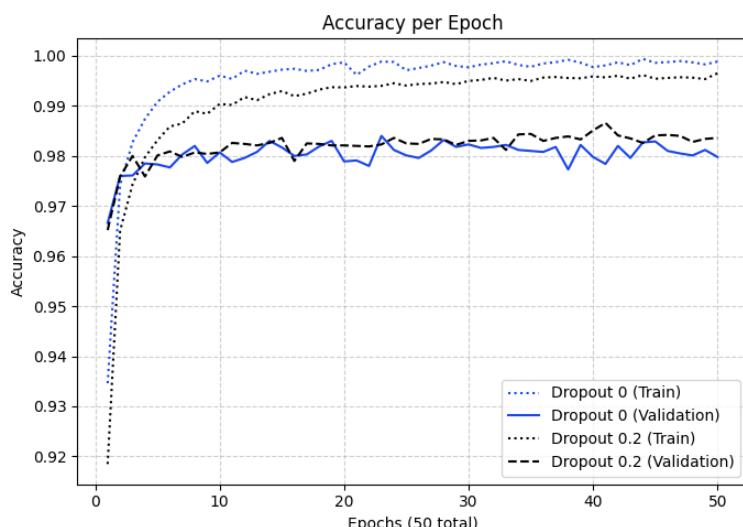
שלוש טכניקות ליבת למניעת Overfitting

Dropout .1

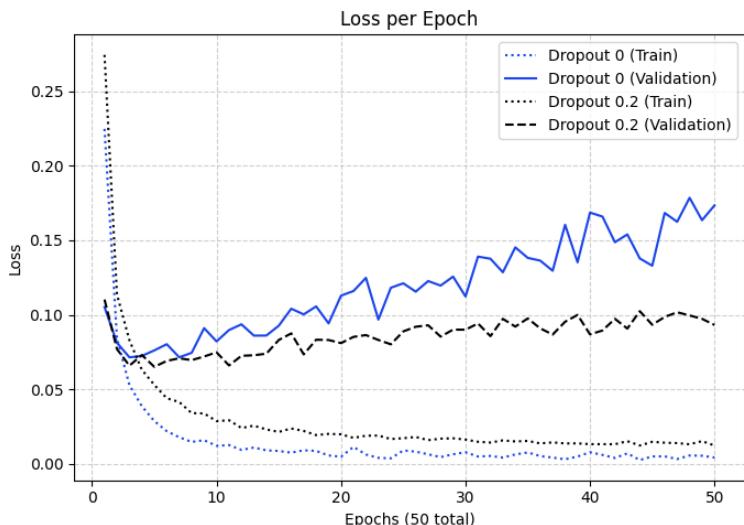
Dropout היא טכניקת רגולרייזציה פשוטה אך עצמאלית. במהלך האימון, כל נeuron חבוי נופל - כלומר מושבת זמנית - בהסתברות p . כך, בכל מעבר קדימה מתאמנותות תתי-רשומות שונות.

פורמלית, אם $h^{(l)}$ הן האקטיבציות בשכבה l , אז Dropout מגדיר מסיכה בינהרית $m^{(l)} \sim \text{Bernoulli}(1 - p)$. בשלב Inference, כל הנוירונים פעילים והמשקלות מותאמות בקנה מידה כדי לפחות על היעדר הדרופ-אואוט.

מנוע הסתגלות-יתר של נוירונים ומעודד פיתוח ייצוגים רובייטיים וחופפים. התוצאה - שיפור בהכללה.



איור 6.1: דיק לאורך Epochs. כחול: ללא דרופ-אואוט. מוקווקו שחורה: דרופ-אואוט עם $p = 0.2$.



איור 6.2: ולידציה Loss לאורך Epochs. כחול: ללא דרופ-אוט. שחור: דרופ-אוט עם $p = 0.2$.

Early Stopping .2

עוקב אחריו Validation Loss בזמן האימון. אם לא נצפה שיפור לאורך מספר מוגדר של אפוקים (נקרא Patience), האימון נעצר מוקדם. אסטרטגיה פשוטה זו מונעת מהמודל להמשיך ולמזרע את Training Loss על חשבון הכללה. היא שימושית במיוחד כאשר אפוקים האופטימלי אינם ידוע מראש.

(L2 Regularization) Weight Decay .3

L2 Regularization, המכונה גם Weight Decay, מענישה משקלות גדולות באמצעות הוספת איבר ריבועי לפונקציית המבחן:

$$(6.1) \quad L_{\text{reg}} = L + \lambda \sum_i w_i^2$$

כאן L הוא loss המקורי (למשל, Cross-Entropy), w_i הוא המשקל

ה-? ו-ג' הוא מקדם הרגולרייזציה. האינטואיציה: משקלות גדולות מצביעות על תלות חזקה בפיצ'רים מסוימים. עניות משקלות כאלה מעודדת את הרשות לפזר את החשיבות בצורה שוויונית יותר ובכך מובילה ללמידה יציבה ומוכללת.

טכנית	מנגנון	מטרה
Dropout	משמעות נוירונים במהלך האימון	מנע התאמת-יתר, מעודד חיפוי ורוביוטיות
Early Stopping	עצר אימון Validation Loss כשי-Loss מפסיק להשתפר	מנע אופטימיזציה עודפת מעבר לנקודת ההקללה
L2 Regularization (Weight Decay)	מוסיף איבר ענייה לפונקציית המבחן $\sum_i w_i^2 \lambda$	מרכז משקלות גדולות, מעודד פתרונות חלקים

טבלה 6.1: השוואת טכניקות רגולרייזציה נפוצות.

מתי לישם רגולרייזציה?

בתרחישים רבים בעולם האמיתי, במיוחד כאשר דאטאsett הולידציה קטן או רועש, קשה לאזות בזדאות Overfitting. לכן, לרוב מומלץ לישם טכניקות רגולרייזציה כברירת מחדל - כאמור מניעה, לא רק כתגובה לכישלון.

סיכום

בשיעור זה בחנו את אחד האתגרים החשובים ביותר בלמידה עמוקה: Overfitting והגדרנו אותו, זיהינו אותו באמצעות עקומות Loss והצגנו שלושה כלים עוצמאים להתמודדות: Weight Decay, Early Stopping ו-Dropout.

טכניקות אלו אינן רק מגבלות את המודל, אלא מנחות אותו ללמידה דפוסים שימושיים במקום לשנן רעש. במודול הבא נעמיק בטכניקות אימון מתקדמות וארכיטקטורות שאפשרות רשתות עמוקות וمبرטות יותר - מבלי לוותר על יכולת ההקללה.

Unit 7

מדוע אנחנו צריכים קומבולוציה?

מבוא

בשיעורים הקודמים בחנו כיצד רשתות נירונים לומדות מדעתה מובנה. בשיעור זה נתמקד בתמונות - תחום עם מבנה מרחבי חזק. רשתות (MLPs) מתעלמות ממבנה זה ודורשות מספר עצום של פרמטרים (משקלות). רשתות קומבולוציה (CNNs) מציאות פתרון יעיל וסקלابلיל יותר. אך מדוע בכלל דרושה קומבולוציה?

שכבות Fully Connected והמגבילות שלן

נבחן תמונה RGB קטנה בגודל 32×32 פיקסלים. מספר הפיצרים בקלט הוא:

$$(7.1) \quad 3 \times 32 \times 32 = 3,072$$

אם וקטור זה משוטח ומחובר לשכבה צפופה של 100 נירונים, מספר המשקלות בשכבה הראשונה בלבד יהיה:

$$(7.2) \quad W_{Total} = 3072 \times 100 = 307,200$$

חיבור מלא זה מאבד מידע מרחבי - הרשת אינה יודעת שפיקסלים סמוכים קשורים זה לזה. ככל שנוסף שכבות נוספות, מספר הפרמטרים

יגדל בmphירות ועלול להוביל ל-Overfitting.

טבלה 1.7: מספר פרמטרים ומודעות מרחבית

מידע מרחביבי	מספר פרמטרים	ארכיטקטורה
✗	307,200	MLP (100 נוירונים)
✓	864	CNN (32 פילטרים של 3×3)

מדוע קונבולוציה עובדת Inductive Bias -

במדעי ה-ML Machine Learning ידוע משפט בשם No Free Lunch הוא קובע שאם מחשבים ממוצע ביצועי אלגוריתם על פני כל הפונקציות האפשריות לצירוף נתונים, אף שיטה אינה טובה מאשר אחרת. הסיבה היחידה לכך ש-CNNs עדיפה על MLPs בתמונות היא ההנחה על **локליות ואינווריאנטיות להזזה**. באופן דומה, ארכיטקטורות אחרות (כמו Transformers) עובדות טוב עם טקסט בזכות היכולת לדגס יחסים בין כל זוג טוקנים.

בפועל, המשמעות היא שאין מודל אחד שהוא "הטוב ביותר" בכל מצב, אלא מודלים שמתאימים למבנה של הדadata שלו. ארכיטקטורות קונבולוציה מצילות במשימות ראייה בדיק מפני שהן מוקדדות הנחות על האופן שבו תבניות ויזואליות מתחנכות.

אופרטור הקונבולוציה

קונבולוציה מפעילה פילטר קטן (שאביריו נלמדים) על אזורים מקומיים בתמונה. נניח $I \in \mathbb{R}^{H \times W}$ היא תמונה קלט בגוני אפור ו- $K \in \mathbb{R}^{k \times k}$ הוא Kernel (פילטר). הקונבולוציה הדיסקרטית הדומיננטית מוגדרת כך:

$$(7.3) \quad F(i, j) = \sum_{u=1}^k \sum_{v=1}^k K(u, v) \cdot I(i+u-1, j+v-1)$$

כאשר:

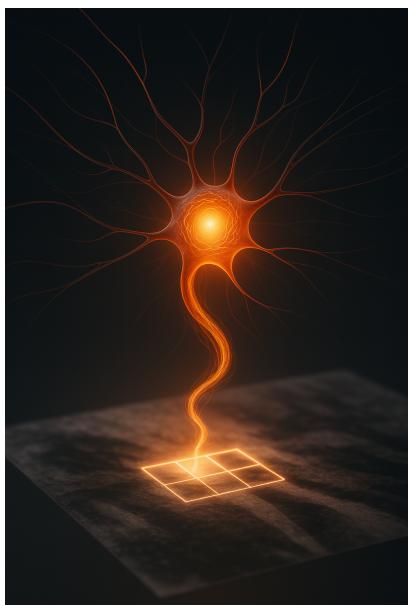
- $F(i, j)$ הוא ערך המיפוי הפיצרי ביציאה במיקום (i, j) .

• $K(u, v)$ הוא משקל הפילטר במיקום (u, v) .

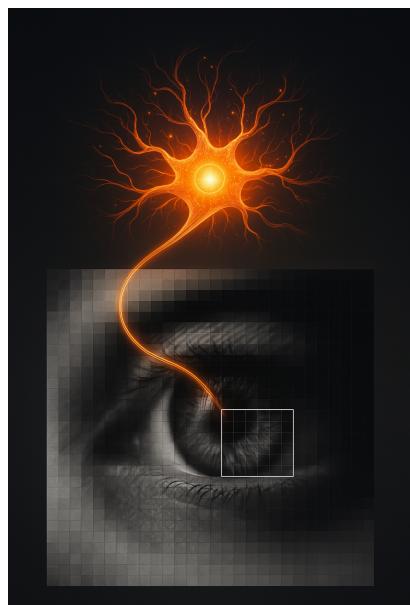
• $I(i + u - 1, j + v - 1)$ הוא ערך הפיקסל בклט עבור חלון הפילטר.

פעולה זו חוזרת על כל המיקומים המרוחביים, תוך שימוש באותו פילטר

- מה שמקטין מאוד את מספר הפרמטרים הנלמדים.



איור 7.2: שימוש חואר בפילטר איננו ריאנטיות להזזה.



איור 7.1: פילטר המופעל על אזור העין - גילי תבנית מקומית.

אינטואיציה של קונבולוציה

קונבולוציה מחלצת תבניות לשימוש חואר מסוורים מקומיים. למשל, פילטר עשוי להיות "עין" או "פינה" שמופיעים בכל מקום בתמונה. בדומה למוח האנושי שאינו משנה אובייקטים שלמים אלא לומד את רכיביהם (קווים, קשתות, עיניים), CNNs מתחמות בזיהוי התבניות שאינן תלויות במיקום מרוחבי.

מאפיינים

אחד היתרונות המרכזיים של קונבולוציה הוא **שיטות פרמטרים** - אותה קבוצת משקولات מופעלת על מיקומים שונים במרחב. כך מספר הפרמטרים הנלמדים קטן מאוד והיעילות החישובית משתפרת. בנוסף, מודלים קונבולוציוניים רוביוטיים מטבעם להזזה - הם יכולים לאזות עצמים ללא תלות במיקום ובתמונה - מה שימושו לשיפור בהכללה ובייציבות. למרות שקונבולוציה עשויה להקטין ממדים מרחבים, התנהוגותה נשלטה גם על ידי:

- **Stride** (א): מספר הפיקסלים שהפילטר זע בכל צעד.
- **Padding** (ב): מספר הפיקסלים המוספים לקצוות כדי לשמור את גודל הפלט.

הגדרות אלו ינתנו פורמלית בשיעור הבא.

CIFAR-10

כדי להשוות בין MLPs ו-CNNs, השתמש בبنץ'מרק **CIFAR-10** - DATASET סטנדרטי לסיוג תמונות. הוא כולל 60,000 תמונות קבוע בגודל $32 \times 32 \times 3$, בעשר מחלקות כגון מטוס, כלב ורכב. הסיבות לשימוש הנפוץ ב-CIFAR-10:

- קטן מספיק לניסויים מהירים.
- כולל שונות אמיתית מעולם האמתי בעשר מחלקות מגוונות.
- CNNs מתקדמים מוגעים לדיקן מעל 90%.



איור 7.3: דוגמאות מותוך. CIFAR-10.

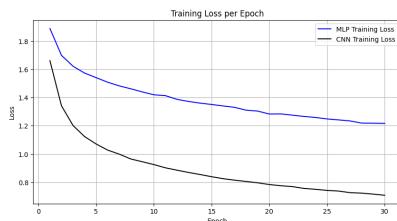
מקרה לדוגמה: MLP לעומת CNN על CIFAR-10

- אימנו שני מודלים על CIFAR-10 במשך 30 אפוקים:
1. MLP מלא עם 2 שכבות צפופות וакטיבציה ReLU.
 2. CNN עם שתי שכבות קונבולוציה וакטיבציה ReLU.
- תוצאות:**

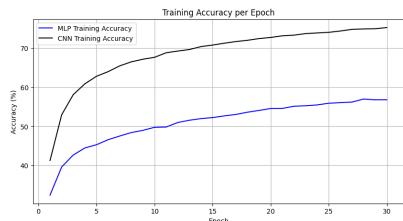
- MLP הגיע רק לדיקוק של ~50%.
- CNN הגיע לכ- 68%, טוב בהרבה.

סיכום

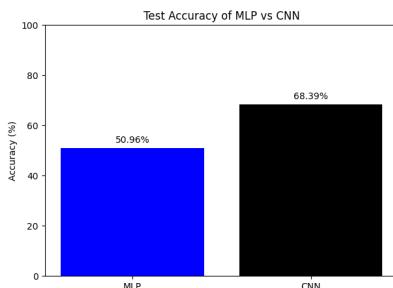
קונבולוציה מציגה Inductive Bias חזק המותאם לדאטה ויזואלי. בניגוד לרשומות MLPs שמתעלמות מהמבנה המרחבי, CNNs שומרות על לוקאליות של פיצ'רים על ידי עיבוד מידע באזוריים מקומיים חופפים. הדבר מאפשר לנו למדוד תבניות מרחביות משמעותיות כגון קצוטות, טקסטורות וצורות. בזכות מנגנון שיתוף המשקولات, CNNs דורשות הרבה פחות פרמטרים -



איור 7.5: מחיר באימון: CNN (שחור) מגע לשגיאה נמוכה יותר מאשר MLP (כחול), מהר יותר.



איור 7.4: דיק באימון: MLP (שחור) מגע לדיק גבוה מאשר CNN (כחול), מהר יותר.



איור 7.6: מחיר באימון - CNN מתכנס מהר יותר ומגע לשגיאה נמוכה יותר.

מה שמאפשר שיפור בהצללה וגם אימון עיל יותר. יתרון נוסף הוא יכולתן לאזהות תבניות בכל מקום בתמונה, מה שמספק רובוטיות להזאה. תוכנות אלו הופכות את CNNs ליעילות וסקלאbilיות מאוד עבור דאטה ויוזאלי מהעולם האמתי, גם כאשר ממדיהם הקלט גדלים.

בשיעור הבא ננסח פורמלית את ההגדרות המתמטיות של Padding, Stride ו-Pooling ונבנה ארכיטקטורות קונבולוציה מלאות למשימות שונות.

שיעור 8

איך פועלת רשת CNN?

מבוא

כעת, לאחר שהבנו מדוע קומבינציה חיונית, נפנה לבנייה הפנימית של רשת קומבינציה (CNN). בשיעור זה נגדיר את רכיבי שכבת הקומבינציה, נסביר כיצד פילטרים גולשים על פני התמונה ונציג מנגנוןים קריטיים כגוןStride, Padding ו-Pooling.

CNNs פועלות על ידי חילוץ פיצ'רים מקומיים באמצעות פילטרים נלדיים. פילטרים אלה מבצעים קומבינציה על פני התמונה ומיצרים Feature Maps - כל אחד מהם מדגיש תבנית אחרת כגון קצוטות, פינות או טקסטורות. כאשר שכבות קומבינציה נערכות באופן היררכי, הן מאפשרות הבנה סמנטית עמוקה של דата חזותי.

שכבת קומבינציה: פילטרים וחלוון הזזה

כל שכבת קומבינציה מורכבת ממספר פילטרים (או Kernels). כל פילטר הוא מטריצה נלמדת של משקولات, בדרך כלל בגודל $k \times k$, למשל 3×3 . בכל מקום מרחבוי, הפילטר מבצע כפל איבר-אייבר עם הקלט המתאים, מסכם את התוצאה ומזין אותה למפת הפיצ'רים ביציאה. פעולה זו נקראת Sliding Window Convolution. פורמלית, נניח $I \in \mathbb{R}^{H \times W}$ היא תמונה גווני אפור ו- $K \in \mathbb{R}^{k \times k}$ הוא

פילטר קונבולוציה. הפלט $F \in \mathbb{R}^{H_o \times W_o}$ מוגדר כך:

$$(8.1) \quad F(i, j) = \sum_{u=1}^k \sum_{v=1}^k K(u, v) \cdot I(i+u-1, j+v-1)$$

כאן (i, j) הוא ערך הפלט במיקום $K(u, v)$, (i, j) הוא מקדם הפילטר במיקום (u, v) ו- $I(i+u-1, j+v-1)$ הוא פיקסל הקלט בחלון הפילטר. הפעלת M פילטרים כאלה מייצרת M Feature Maps - אחד לכל תבנית שהרשת למדעה אלה.

והשפעות קצה Padding

ללא Padding, הפילטר אינו יכול לכסות את אזורי הקצה של התמונה במלואם. כתוצאה לכך, כל קונבולוציה מקטינה את הממדים המרחביים של מפת הפייצ'רים. כדי לשמר את הממדים, מושגים Zero-Padding - הוספה שורות ועמודות של אפסים סביב הקלט.

$$x = \begin{bmatrix} ? & ? \\ 1 & 2 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix} \rightarrow \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 2 & 3 & 4 & 0 \\ 0 & 5 & 6 & 7 & 8 & 0 \\ 0 & 9 & 10 & 11 & 12 & 0 \\ 0 & 13 & 14 & 15 & 16 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

איור 8.1 : קצוטות נוספים מאפשרים CISIOI מלא של מימדי הקלט.

Padding חשוב במיוחד במקרים עמוקות, שבו שכבות קונבולוציה רבות עלולות להקטין את הרזולוציה המרחבית יתר על המידה.

והשפעה Stride

ה-Stride s קבוע בכמה פיקסלים הפילטר יז בכל צעדי. ערך 1 אומר שהפילטר יז פיקסל אחד בכל פעם - והוא פלט ברזולוציה גבוהה. ערך 2 מدلג על

מיקומים חלופיים, מצמצם את ממד הפלט ומייעל את החישוב.

$$\text{Stride = 1} \quad x = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix} \quad x = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix} \quad x = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix}$$

$$\text{Stride = 2} \quad x = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix} \quad x = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix}$$

- איור 8.2: השפעת השפה Stride=1 - למעלה (דגימה צפופה), למטה (דגימה גסה יותר).

כעת, בשילוב בין גובה הקלט H , p Padding, גודל הפילטר k וה-Stride s , ממד הפלט מוגדרים כך:

$$(8.2) \quad H_{\text{out}} = \left\lfloor \frac{H + 2p - k}{s} + 1 \right\rfloor$$

$$(8.3) \quad W_{\text{out}} = \left\lfloor \frac{W + 2p - k}{s} + 1 \right\rfloor$$

Pooling

לאחר קונבולוציה, לעיתים קרובות מבצעים פעולה של Pooling כדי לדוגם מחדש את מפת הפייצרים. הזרה הנפוצה ביותר היא **Max Pooling** - בחרית הערך הגבוה ביותר מכל טלאי מקומי. לחלוfin, **Average Pooling** מחשב את הממוצע באזור.

פורמלית, אם $R_{i,j}$ הוא אזור pooling עבור מיקום הפלט (i, j) , אז:

:Max Pooling

$$(8.4) \quad P_{\max}(i, j) = \max_{(u,v) \in R_{i,j}} F(u, v)$$

:Average Pooling

$$(8.5) \quad P_{\text{avg}}(i, j) = \frac{1}{|R_{i,j}|} \sum_{(u,v) \in R_{i,j}} F(u, v)$$

$$x = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \rightarrow 4$$

$$x = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \rightarrow 2.5$$

.איור 8.3: גרסאות Max Pooling למעלה , למיטה Pooling . לשתייהן מקטיניות רזולוציה מרחבית.

בלוק CNN בסיסי

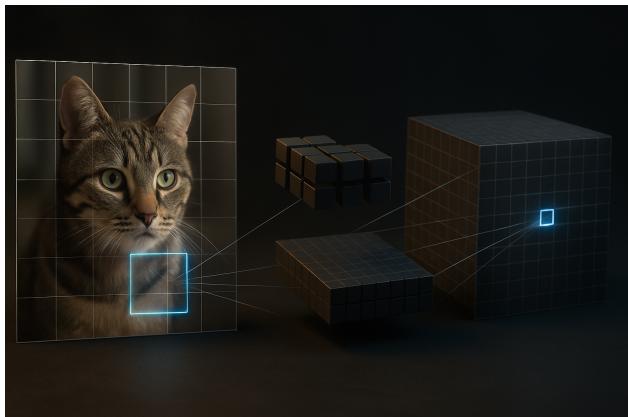
בלוק CNN טיפוסי כולל את השלבים הבאים:

- **קלט:** תמונה, למשל בגודל 28×28 פיקסלים.
- **Convolution :** מימוש פילטרים 3×3 (למשל 8 פילטרים נלמדים).
- **ReLU :** יישום פונקציה לא-ליניארית, לרובה Activation .
- **Pooling :** הקטנת הממדים המרחביים באמצעות Max Pooling .
- **Flatten :** הרמת מפת הפיצ'רים לווקטור חד-ממדי.
- **Fully Connected Layer :** ביצוע סיווג סופי על בסיס הפיצ'רים שנחולצו.

בלוקים כאלה מרכיבים יחד מודלים עמוקים יותר. עם כל שכבה נוספת, גודל Receptive Field של הפלט גדל

עומק Receptive Field

למרותSCP של יחידת קונבולוציה רואה רק אזור קטן, אוסף שכבות רבות מגדילה את ה-*Receptive Field* - האזור בקלט שמשפיע על נוירון ביציאה.



איור 8.4: ה- Receptive Field גדול עם עומק - מה שמאפשר למודל ללמידה גלובלי.

ברשת CNN השכבות התחתנות מזהות קצויות מקומיים, שכבות ביניות מזהות טקסטורות ומוטיבים והשכבות העליונות מגיבות לאובייקטים שלמים.

סיכום

CNN משלבת פילטרים, Padding ו-Stride Pooling למערכת היררכית חזקה ללמידה פיצ'רים מרחביים. מנגנונים אלה מאפשרים לרשת להקטין את גודל הדאטה, לשמר מבנה ולבנות ייצוגים מופשטים יותר ויוטר של קלט חזותי. בשיעור הבא נציג שיפורים ארכיטקטוניים כגון Batch Normalization ו-Dropout המשפרים עוד יותר את הביצועים והrobustיות.

শיעור 9

רוצפים וזמן: LSTM , GRU ו-RNN

מבוא

עד כה עסקנו בקלטים כגון תמונות, שבהם כל הדadata זמין בבת אחת. אך תחומיים רבים - כגון שפה טبيعית, מזיקה וזרמי חישנים - דורשים עיבוד **דאטה רציף** (Sequential Data). במשימות כאלה ההקשר הtemporal חיוני: המודל חייב לא רק "לראות" אלא גם **לזכור**.

בשיעור זה נציג את משפחת רשותות RNNs (Recurrent Neural Networks), נתחיל מה-RNN הבסיסי, נתხך את מוגבלותיו ואז נציג שתי ארכיטקטורות נפוצות: Bi-Directional LSTM ו-GRU.

RNN רשותות

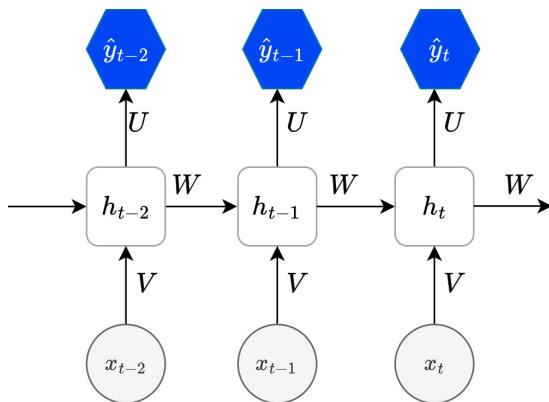
בניגוד לתמונות, שבחן כל הפיקסלים מעובדים בו-זמןית, רוצפים נפרשים צעד אחר צעד. בכל צעד זמן t , נכנס קלט חדש x_t והמודול חייב לעדכן את מצבו הפנימי כך שיישקף גם את הקלט הנוכחי וגם את ההיסטוריה הקודמת. RNN שומרת מצב חבי h_t שמתעדכן בצורה רקורסיבית:

$$(9.1) \quad h_t = f_h(\mathbf{V}x_t + \mathbf{W}h_{t-1} + b_h)$$

$$(9.2) \quad \hat{y}_t = f_y(\mathbf{U}h_t + b_y)$$

כאשר:

- x_t הוא הקלט בצעד זמן t .
- h_t הוא המצב החבוי בצעד זמן t .
- $\mathbf{U}, \mathbf{V}, \mathbf{W}$ הן מטריצות משקלים נלמדות.
- b_h, b_y הן הטוויות.
- f_h היא פונקציית אקטיבציה (למשל tanh או ReLU).
- \hat{y}_t הוא הפלט של המודל בצעד t .



איור 9.1: ארכיטקטורת RNN לאורך זמן.

מנגנון זה מאפשר למידה לאורום דרך הזמן. עם זאת, RNNs קלאסיות סובולות מבעית **Vanishing Gradient**, במיוחד על רצפים ארוכים. הגרדיינטים דועכים בעת Backpropagation Through Time (BPTT), מה שמניביל את יכולת הלמידה תלות ארוכת טווח.

LSTM: Long Short-Term Memory

רשת LSTM מציגה תא זיכרון פנימי C_t ומערכת שערים השולטים במעבר המידע. ארכיטקטורה זו נבנתה כדי להתמודד עם

בעיית ה- Vanishing Gradient מוגדרים כך:

$$(9.3) \quad i_t = \sigma(W_i x_t + U_i h_{t-1} + b_i) \quad (\text{שער קלט})$$

$$(9.4) \quad f_t = \sigma(W_f x_t + U_f h_{t-1} + b_f) \quad (\text{שער שכחה})$$

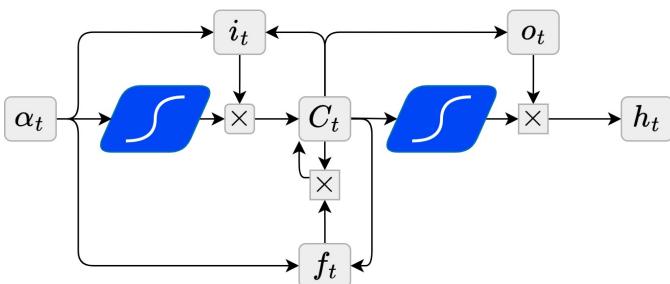
$$(9.5) \quad o_t = \sigma(W_o x_t + U_o h_{t-1} + b_o) \quad (\text{שער פלט})$$

$$(9.6) \quad \tilde{C}_t = \tanh(W_c x_t + U_c h_{t-1} + b_c)$$

$$(9.7) \quad C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t$$

$$(9.8) \quad h_t = o_t \odot \tanh(C_t)$$

כאן σ היא פונקציית סיגמוואיד, \odot מסמל כפל איבר-איבר, C_t הוא מצב הטה ו- h_t הוא המצב החבוי.



איור 9.2: המבנה הפנימי של יחידת LSTM.

ארQUITקטורה זו מאפשרת זכרון ושכחה סלקטיבית והופכת את ה-LSTM לאידיאלי ללמידה תלות קצרה וארוכה אחד.

GRU:Gated Recurrent Unit

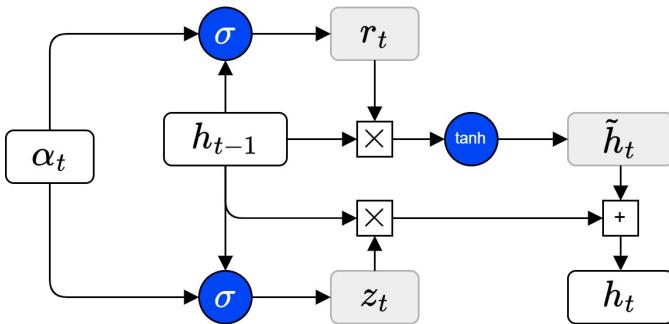
GRU היא גרסה פשוטה יותר של LSTM המזגת את שער הקלט ושער השכחה לשער עדכון יחיד. יש לה פחות פרמטרים ולעיטים מתאימים מהר יותר, תוך שמירה על ביצועים דומים. משוואות ה-GRU הן:

$$(9.9) \quad z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z) \quad (\text{שער עדכון})$$

$$(9.10) \quad r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r) \quad (\text{שער איפוס})$$

$$(9.11) \quad \tilde{h}_t = \tanh(W_h x_t + U_h(r_t \odot h_{t-1}) + b_h)$$

$$(9.12) \quad h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$



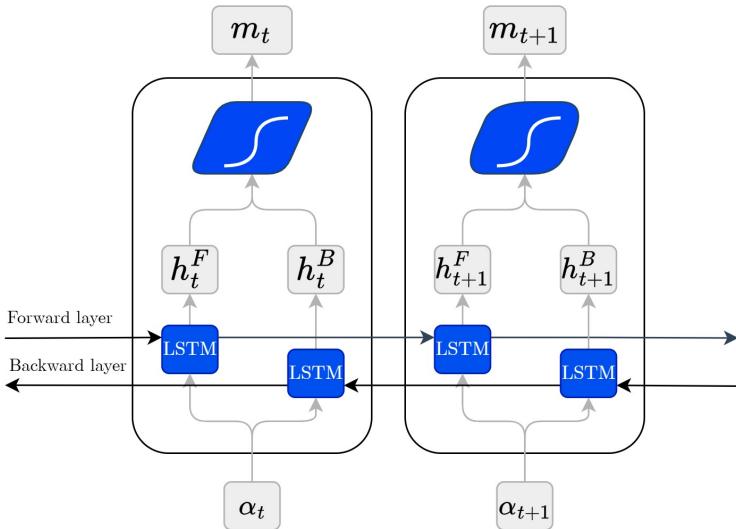
איור 9.3: ארכיטקטורת GRU עם שני שערים: עדכון ואיפוס.

זולות חישובית ולעיטים מועדות עבור רצפים ארוכים כאשר זמן האימון מהוות מוגבלת.

Bi-Directional LSTM

חלק מההשימות מutowichtet מכישה גם לעבר וגם לעתיד. Bi-Directional LSTM (Bi-LSTM) משיגה זאת על ידי הרצת שתי LSTMs במקביל: אחת מעבדת קדימה ($t = 1 \rightarrow T$) והשנייה אחורה ($t = T \rightarrow 1$). הפלטים משולבים בכל צעד זמן:

$$(9.13) \quad h_t^{\text{bi}} = [h_t^{\rightarrow}; h_t^{\leftarrow}]$$



איור 9.4: Bi-LSTM : שילוב מצבים חבויים קדימה ואחורה.

Bi-LSTMs ייעילות במיוחד במשימות עיבוד טקסט, שבן משמעות המילה תלויה לעיתים גם במילים שקדמו וגם באלה שבאו אחריה.

דוגמה: ניתוח סנטימנט

נבחן את המשפט:

“The service was okay, but the food was absolutely terrible.”

RNN חד-כיווני המעבד משמאל לימין עשוי לקרוא תחילת את החלק *“The service was okay”* ולהסיק שהסנטימנט ניטרלי או חיובי כמעט. אך האות המרכזית מופיעה רק בסוף: *“absolutely terrible”*, המעיד על שליליות חזקה.

RNN חד-כיווני עלול להיכשל בתיקון התחזית המוקדמת. לעומת זאת, Bi-LSTM מעבדת את הרצף בשני הכוונים. היא יכולה לשלב את ההקשר

מהסוף ("e") בהבנתה את המשפט כולו - וכך להגיע לסיוג סנטימנט מדויק יותר. דוגמה זו מדגימה את התלות בהקשר.

מודל	shoreuds	זיכרון	יתרונות	שימוש
RNN LSTM GRU	אין קלט, שכחה, פלט	קצר טווח ארוך טווח	פשט יציב, חזק ארוך טווח	רצפים קצרים טקסט, אודיו כמו LSTM

טבלה 9.1: סיכום ארכיטקטורות למידע רציף.

סיכום

מודלים רציפים חיוניים למשימות הכוללות דатаה בזמן. ה-RNN הבסיסי מציג את עקרון הזיכרון לאורך זמן. LSTM ו-GRU משפרים זאת על ידי מתן יכולת לשמר ולשלוט במצב בצורה עיליה יותר.

במשימות הדורשות הקשר עתידי, Bi-LSTM מרחיבה יכולת זו. עם זאת, ארכיטקטורות אלו עדין נשענות על חישוב סדרתי, מה שמנבל את הסקלאbilיות שלהם.

בשיעורים הבאים נציג את ארכיטקטורת ה-Transformer ששולטת כיום בשיטות המתקדמות ביותר לעיבוד שפה ומשימות רצף.

שיעור 10

הורדת מימדים עם Autoencoders

מבוא

עד כה התמקדנו בלמידה מונחית (Supervised Learning) - אימון מודלים על זוגות קלט-פלט עם תוצאות. בשיעור זה נחקור פרדיגמה אחרת: **למידה לא מונחית** (Unsupervised Learning), שבה אין תוצאות נתנות והמטרה היא לחשוף מבנה חבוי בתוך הדadata.

באופן כללי, למידת מכונה מחולקת לשלווה סוגים מרכזיים:
1. המודל מפנה קלטים x לפלאטים y ; *Supervised Learning* (1)
2. מציאת דפוסים או דחיסה ללא תוצאות; *Unsupervised Learning* (2)
3. סוכן הלומד דרך אינטראקציה ותגמול. *Reinforcement Learning* (3)

אמנם הקורס מתמקד בלמידה מונחית, אך מודלים לא-מנחים חיוניים באותה מידה. ביןיהם, (AE) **Autoencoder** ממלא תפקיד מרכזי בלמידה عمוקה מודרנית. הוא רשות נוירונים המתאמת לשחזור את הקלט המקורי - ובתהליך זה היא לומדת ייצוג פנימי של הדadata.

ארכיטקטורת Autoencoder

מודול Autoencoder מרכיב משתי תת-rstות:

Encoder • f_{enc} : דוחס את הקלט $x \in \mathbb{R}^d$ לווקטור לטני $z \in \mathbb{R}^k$, כאשר $k \ll d$

• Decoder: משחזר את הקלט $\hat{x} \in \mathbb{R}^d$ מתוך z .

פורמלית, הפונקציות מקיימות:

$$(10.1) \quad z = f_{\text{enc}}(x)$$

$$(10.2) \quad \hat{x} = f_{\text{dec}}(z)$$

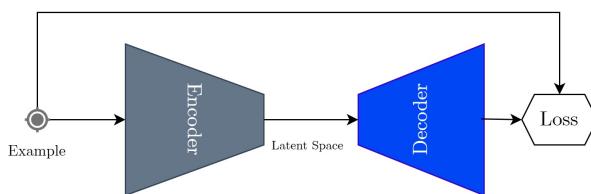
הרשת מתאמת לפחות את שגיאת השחזור:

$$(10.3) \quad \mathcal{L}(x, \hat{x}) = \|x - \hat{x}\|^2 = \sum_{i=1}^d (x_i - \hat{x}_i)^2$$

פונקציית מחיר זו מעודדת את המודל לשמור מידע חיוני ב- z , תוך השלcta
מידע לא רלוונטי.

למידת "צוגים קומפקטיים"

על ידי שימוש במבנה של צוואר בקבוק (כלומר $d \ll k$), ה-Autoencoder נאלץ ללמידה רק את הדפוסים המשמעותיים ביותר בקלט. ה-Encoder מחלץ פיצ'רים מופשטים וה-Decoder לומד להרכיב אותם מחדש מודיעק של הקלט המקורי. הדבר יעל במיוחד עבור נתונים עתיר ממדים כמו תמונה, אודיו ואותות חישניים, שבהם דחיסה ויכולת הפרוש חשובות.



. איור 10.1: המחתת Autoencoder.

דוגמה: הורדת ממדים ב- MNIST

להדגמת הרעיון, ניישם Autoencoder על Dataset MNIST - המורכב מתמונות גווני אפור בגודל $28 \times 28 = 784$ פיקסלים של ספרות בכתב יד. בניית Autoencoder :

- קלט: $x \in \mathbb{R}^{784}$

- מרחב לטנטי: $z \in \mathbb{R}^2$

- פלט: $\hat{x} \in \mathbb{R}^{784}$

מספר הפרמטרים הנלמדים בראש זו הוא 218,514

Layer (type)	Output Shape	Param #
input_layer (InputLayer)	(None, 784)	0
dense (Dense)	(None, 128)	100,480
dense_1 (Dense)	(None, 64)	8,256
dense_2 (Dense)	(None, 2)	130
dense_3 (Dense)	(None, 64)	192
dense_4 (Dense)	(None, 128)	8,320
dense_5 (Dense)	(None, 784)	101,136

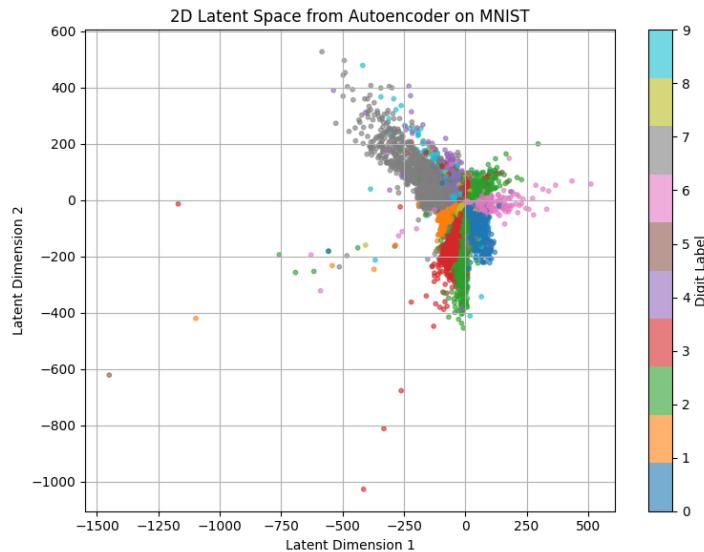
.איור 10.2: ארכיטקטורת Autoencoder עבור MNIST.

לאחר האימון, אנו מחלכים את הווקטורים הלטנטיים z עברו כל הספרות. למרות שלא ניתנו תוויות במהלך האימון, ספרות עם מבנה דומה (למשל "1" ו-"7") מקובצות יחד.

דבר זה מדגים שהמודל הצליח למדוד את המבנה הפנימי של הדאטה - סימן מובהק לomidת ייצוג מוצלח.

פתרונות השחזר

נציג שחזורי ספרות מتوزע ה-Autoencoder. לכל תמונה קלט x , נעביר אותה ב- Encoder-Decoder ונשווה לשחזר \hat{x} : על אף צוואר הבקבוק, המודלצליח לשמור את רוב המבנה, מה שמראה כי חילץ פיצ'רים לטנטאים ממשמעותיים.



איור 10.3: מרחב לטנטי דו-dimensional שנלמד על ידי Autoencoder. הצבעים מציננים תוויות אמיתיות (לצורך ויזואלייזציה בלבד).

עקומת האימון

האיור הבא מציג את שגיאת השחזר (MSE) כפונקציה של מספר האפוקים. ככל שהאימון נמשך, השגיאה יורדת בהדרגה:

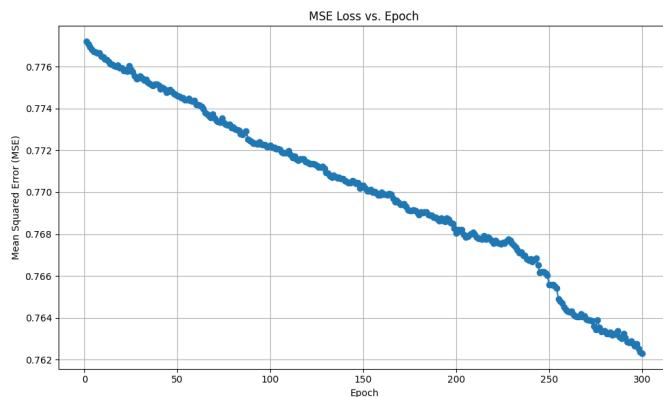
השוואה ל- PCA

PCA (Principal Component Analysis) היא טכניקה LINEARIT למצטום ממדים. היא מוצאת ציוונים אורותוגונליים הממקסימים וריאנס. עם זאת, PCA אינה יכולה למוד א-LINEARITY בDATA.

מכילים את PCA באמצעות פונקציות לא-LINEARITY (למשל AkTIVITIET ReLU), מה שמאפשר להם ללמידה ירידות מעוקلات להפרדת המחלקות למשל. ההבדל המרכזי הוא ש-PCA מוצאת ת-מרחב LINEAR MINIMALI בשגיאת שחזר ריבועית, בעוד ש-Autoencoders לומדים פונקציה לא-LINEARIT המזערת אותה פונקציה. גמישים יותר ומתאים



איור 10.4: ספרות מקוריות (למעלה) ושוחזר (למטה) לאחר אימון Autoencoder ל-300 אפוקים.



איור 10.5: שגיאת שחזור MSE לעומת אפוקים. המודל מתכנס ביציבות לשגיאה נמוכה.

במיוחד לדאטה מובנה כמו תמונות.

יישומים והרחבות

Autoencoders שימושו בהצלחה במגוון רחב של משימות למידה עמוקה. כמו מהיישומים המרכזיים כוללים:

- למידת ייצוגים במרחב נמוך מקלטים במרחב גבוה.
- דחיסת אוטות לאחסון או שידור עיל.
- זיהוי חריגות/ anomalיות בהתבסס על שגיאת שחזור גבוהה.
- פרא-טריינינג לרשותות עמוקות ולאחר מכן Fine-Tuning מונחה.

- ניקוי רעש מתקלט (Denoising).

מעבר ל-⁻Autoencoder הבסיסי, פותחו הרחבות רבות לצרכים ספציפיים. אחת מהן היא Variational Autoencoder (VAE), שבה ה-⁻Encoder מפיק פרמטרים של התפלגות גאוסית - וקטור מומצאים $(x) \mu$ וקטור סטיות תקן $(x) \sigma$ - שמננו נדגם המשטנה הlatentי z כך:

$$(10.4) \quad z \sim \mathcal{N}(\mu(x), \sigma^2(x)).$$

סיכום

Autoencoders הם כלי יסודי בלמידה לא-דוחנחת, החושף את המבנה הפנימי של נתונים באמצעות דחיסה ושחזר. הם אינם נשענים על תוויות ובכל זאת לומדים ייצוגים שימושתיים שניית להמחיש, לקבץ ולהשתמש בהם במשימות המשך.

הדוגמה על MNIST ממחישה כיצד Autoencoders חושפים תכונות גיאומטריות במרחב הספרות מבל שנקמר להם מהי ספרה. כוחם טמון ביכולת להקליל וללמוד את מהות המבנה של הדата. בכך הם הופכים לכלי בלתי נפרד במערכות למידה עמוקה מודרניות.

שיעור 11

Transformer וה-Self-Attention

מבוא

ארQUITטורת ה-Transformer, שהוצגה בשנת 2017 במאמר You Need Attention Is All You Need, חוללה מהפכה בלמידה עמוקה עבור דאטה רציף. במקור פותחה עבור (NLP) Natural Language Processing, אך בהירה הפכה ליסוד של מודלים גנרטיביים מודרניים כגון GPT ו-BERT. החידוש המרכזי שלה - Self-Attention - החליף את הרקורסיה במנגנון המאפשר לכל טוקן להתייחס לכל שאר הטוקנים, מה שאפשר מיקבול רחב יותר ומידול טוב יותר של תלות לטווח טוקן. מאוחר יותר הורחב עקרון זה גם לתחומים אחרים: למשל, Vision Transformers (ViT) על תМОנות.

מדוע RNN לא מספיקות

רשתות RNN וגרסאותיהן כמו LSTM מעבדות רצפים טוקן אחרי טוקן, תוך עדכון המצב החבוי השומר מידע מה עבר. גישה סדרתית זו מקשה על מיקבול ומגבילה את יכולת ללמידה תלות רחוכה.

דוגמה: “*The book that the professor recommended at the end of the lecture was truly brilliant.*” כדי לפרש נכון את המילה “book” כמתיחסת ל-“book”, על המודל ללמידה תלות ארכוכת טווח - משחו ש-RNN מתקשה לשמר.

מנגנון ה-Attention

פותר זאת על ידי כך שכל טוקן יכול להתייחס ישירות לכל הטוקנים האחרים ברכף. נניח $X = [x_1, x_2, \dots, x_T] \in \mathbb{R}^d$, כאשר כל $x_i \in \mathbb{R}^d$. Embedding של הטוקן i -י הוא וקטור קלט מועבר באמצעות ליניארי לשולושה ייצוגים:

$$(11.1) \quad q_i = W_Q x_i$$

$$(11.2) \quad k_i = W_K x_i$$

$$(11.3) \quad v_i = W_V x_i$$

כאשר:

$q_i \in \mathbb{R}^{d_k}$ הוא וקטור ה-Query של הטוקן i .

$k_i \in \mathbb{R}^{d_k}$ הוא וקטור ה-Key של הטוקן i .

$v_i \in \mathbb{R}^{d_v}$ הוא וקטור ה-Value של הטוקן i .

$W_Q, W_K, W_V \in \mathbb{R}^{d_k \times d}$ הן מטריצות נלמדות.

ציוו ה-Attention בין טוקן i לטוקן j מחושב באמצעות מכפלה מנורמלת:

$$(11.4) \quad \alpha_{ij} = \frac{q_i^\top k_j}{\sqrt{d_k}}$$

לאחר מכון מחושבים משקלים ה-Attention באמצעות פונקציית Softmax:

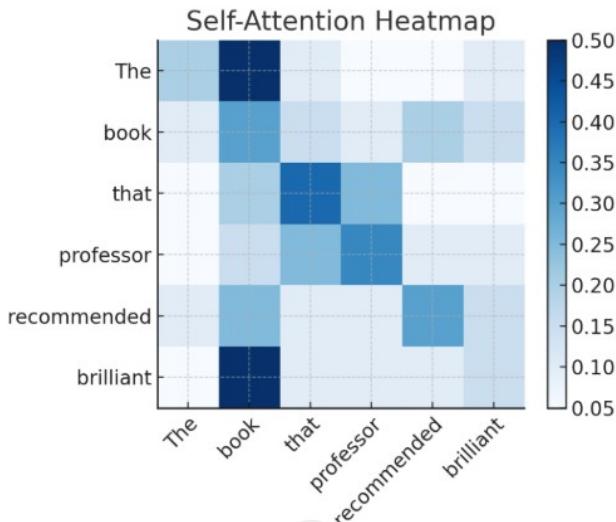
$$(11.5) \quad a_{ij} = \frac{\exp(\alpha_{ij})}{\sum_{j'=1}^T \exp(\alpha_{ij'})}$$

לבסוף, כל ערך מתעדכן על ידי סכום משוקלל של וקטורי ה-Value:

$$(11.6) \quad z_i = \sum_{j=1}^T a_{ij} v_j$$

תהליך זה יוצר ייצוג מותאם-הקשר \tilde{z} עבור כל טוקן.

نبיט לדוגמה במאט Self-Attention עבור המשפט: *The solution that the professor recommended was brilliant.*”
 כל שורה מייצגת את ה-Query של מילה והעמודות מייצגות את המילים שאליהן היא מתייחסת. תאים כהים יותר מצביעים על משקל Attention גבוה יותר. לדוגמה, המילה “book” מתייחסת בחזקה ל-“brilliant” וכן נשמרת ההתאמה למרות טוקנים שישנים רבים המופיעים ביניהם.

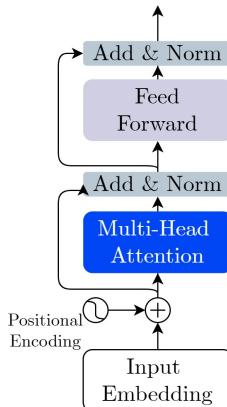


אייר 11.1: מאט Self-Attention עבור המשפט: *The solution that the professor recommended was brilliant.*”

מה-Attention ל-Transformer

Transformer מלא מרכיב **מ-ארכיטקטורת Encoder-Decoder**. ה-Encoder ממפה רצף קלט (למשל משפט) לרצף “יצוגים חבויים באמצעות שכבות של Masked Self-Attention ורשתות Feedforward. ה-Decoder משתמש ב-“Cross-Attention (כדי למנוע הצגה קדימה), ב-“Feedforward עם פלט ה-Encoder ובראשת Feedforward כדי ליצור את פלטי הטוקנים אחד-אחד.

מבנה זה מאפשר הפקה מותנית, מה שהופך את ה-Transformer למודול ג'נרטיבי עצמאי - מתאים למשימות כגון תרגום, יצירתי טקסט, כתיבת תיאורי תכונות ועוד.



איור 11.2: המחשפת Encoder

האיור מציג את המבנה הפנימי של בлок Encoder. בלבו נמצא מודול Multi-Head Attention - הרחבה מקבילית של מנגנון ה-Attention - המאפשר למודל להתייחס לתת-מרחבים שונים של ייצוג הקלט בו-זמנית. פלט ה-Attention עובר דרך Residual Connection, נרחב עליו בהמשך, ואחריו Layer Normalization (המכונה "Add & Norm"), המייצבת את התפלגות האקטיבציות ומאפשרת ארכיטקטורות עמוקות יותר. לאחר שכבת ה-Attention באה רשות Position-Wise Feedforward, בדרך כלל שתי שכבות מלאות עם פונקציית אקטיבציה לא-ליניארית (כגון ReLU או GELU). חיבור שאריות נוספת. נספּ ושלב נורמליזציה משלימים את הבלוק.

חשוב לציין שארכיטקטורה זו מאפשרת מידול הקשר תלוין-סדר מלא ללא שימוש ברקורסיה או קונבולוציה, תוך הסתמכות בלבד על דפוסי Attention נלמדים ועל Positional Encodings. מבנה זה מהווה את עמוד השדרה של מודלים שפתיים מודרניים.

יתרונות והכללות

ל-*Transformers* יתרונות מרכזיים על פני ארכיטקטורות קודמות: הם מאפשרים מיקבול מלא של כל רכיבי הרצף - מה שמאיצץ אימון. מנגנוני ה-*RNNs* ממדלים באופן טבעי הקשר גלובלי ותלות אורך טווח, בניגוד ל-*Pretrained Language Models* (Deep Networks) ותומכים בהעברת המתקשים בכך. הם ניתנים להעמקה רבה (Deep Networks) ותומכים בהעברת ידע דרך מודלים מאומנים מראש (Pretrained Language Models).

סיכום

Self-Attention שינתה מהיסוד את ארכיטקטורת המודלים בלמידה עמוקה. באמצעות מתן אפשרות לכל טוקן לאסוף מידע מכל הרצף, נוצרו מערכות חזקיות וסקלאbilיות שתתרמו לתחומיים רבים. ארכיטקטורת ה-*Transformer*, עם מבנה ה-*Encoder-Decoder*, מהוות את הבסיס למודלים ה'גנרטיביים' המוביילים כיום בתחום עיבוד השפה הטבעית ואך מעבר.

ו שיעור 12

Initialization ו Normalization

מבוא

בשיעור זה נחקרות שתי שיטות קריטיות בלמידה عمוקה: **Normalization** ו-**Initialization**. טכניקות אלו חיוניות לייצוב והאצת האימון של רשתות נוירוניים عمוקות. ללא נורמליזציה ואתחול נוכנים, מודלים עלולים להתכנס לאט, להיתקע בנקודות מינימום מקומיות או לסבול מבעיות של Vanishing Gradients או Exploding Gradients. נתחיל בהבנת הצורך לנורמל אקטיבציות ברשנות עמוקות ולאחר מכן נראה כיצד בחירת משקלות התחלה מושפעת על מעבר המידע בראשת.

מדוע לנורמל?

במהלך האימון, האקטיבציות זורמות דרך הרשת. עם זאת, ככל שהעומק גדל, ההתפלגות שלן עשויה להשתנות בצורה דרמטית בין שכבות או מיני-באצ'ים. תופעה זו, הידועה בשם Internal Covariate Shift, מפריעה לתהליכי הלמידה וגורמת לתנודות בגרדיינטס. כדי למתן זאת, משתמשים בטכניות נורמליזציה כגון **Layer Normalization** ו-**Batch Normalization**.

Batch Normalization

בහינתן מינימליזציה של אקטיבציות $\{x_1, x_2, \dots, x_m\}$ מנירון מסוים (או ערוֹץ פיז'ר), מNORMAL את הערכים על פני ממד הbaru:

$$(12.1) \quad \mu_B = \frac{1}{m} \sum_{i=1}^m x_i$$

$$(12.2) \quad \sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$$

$$(12.3) \quad \hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \varepsilon}}$$

$$(12.4) \quad y_i = \gamma \hat{x}_i + \beta$$

כאשר:

- μ_B ו- σ_B^2 הם הממוצע והשונות המוחسبים על פני המינימליזציה. \hat{x}_i היא האקטיבציה המנורמלת.
- γ ו- β הם פרמטרים נלמדים המאפשרים שינוי קנה מידת והזזה. ε הוא קבוע קטן להבטחת יציבות נומרית.

בדרך כלל מוסיפים BatchNorm אחרי הטרנספורמציה הלייניארית ולפני פונקציית האקטיבציה:

Linear \rightarrow BatchNorm \rightarrow ReLU

שיטת זו משפרת את חישוב הגרדיאנטס, מאפשרת אימון מהיר יותר ואף מאפשרת שימוש בקצב למידה גבוה יותר. עם זאת, היא משתמשת על מינימליזציות מספיק גדולים לסטטיסטיקה יציבה.

Layer Normalization

במשימות עם מיני-batch'ים קטנים מאוד או ארכיטקטורות מבוססות רצף (כמו מודלים לשפה טבעיות או Transformers), BatchNorm עלול להיות פחותiesel בשל סטטיסטיות באז' לא יציבות. במקרה אלון, **Layer Normalization** מתאים יותר, משום שהוא מנормל על פני הפיצ'רים של כל דוגמה בודדת - ולא על פני הבאז'.

בהינתן וקטור קלט $x \in \mathbb{R}^d$ (למשל האקטיבציות של שכבה אחת עבור דוגמה בודדת), החישוב הוא:

$$(12.5) \quad \mu = \frac{1}{d} \sum_{i=1}^d x_i$$

$$(12.6) \quad \sigma^2 = \frac{1}{d} \sum_{i=1}^d (x_i - \mu)^2$$

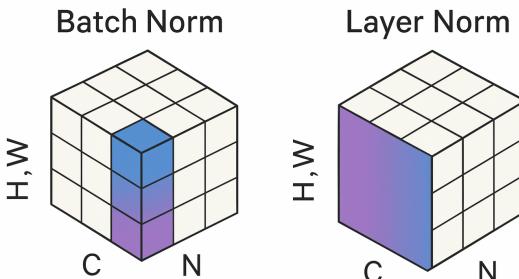
$$(12.7) \quad \hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$(12.8) \quad y_i = \gamma \hat{x}_i + \beta$$

כאן, הנורמליזציה מבוצעת על פני כל d הפיצ'רים של הדוגמה היחידה ואוותם פרמטרי קנה מידה והזה נלמדים (β, γ). שיטה זו אינה תלואה בגודל הבאז' ועובדת היטב במקרים שבחם נדרשת עקבות לכל דוגמה, כמו ברשתות רצף (RNNs) או דקודרים אוטורגרסיביים. היא אפקטיבית במיוחד ב-Transformers וברשתות רצף שבהן מבניות הקלט משתנות ותלות בזמן.

Weight Initialization

לפני תחילת האימון, המשקלות מאותחלות לרוב בערכים אקראים. עם זאת, אתחול לא נכוון עלול לעוזת את זרימות המידע, לגורום ל-Vanishing/Exploding Gradients ולעכוב את הלמידה. המטריה היא לשמר שונות מבוקרת של



אייר 12.1: השוואה בין Norm Layer (משמאל) ו-Batch Norm (מימין). Norm מנורמל על פניו מממד הבאץ' (N) בעוד Norm Layer מנורמל על פניו כל הפיצ'רים בדוגמה אחת (C, H, W).

הакטיבציות והגרדיינטיהם בכל השכבות.

Xavier Initialization

שיטת זו מתאימה לאקטיבציות סימטריות כמו \tanh או sigmoid . דוגמת משקלות מהתפלגות עם שונות:

$$(12.9) \quad W_{ij} \sim \mathcal{N}\left(0, \frac{2}{n_{\text{in}} + n_{\text{out}}}\right)$$

כאשר n_{in} ו- n_{out} הם מספר הנוירונים הנכנסים והויצאים בשכבה. שיטה זו מאזנת את מעבר המידע קדימה ואחוריה ברשות כדי למנוע דעיכה או הגדלה יתר על המידה של הנתונים.

He Initialization

עבור אקטיבציות לא-סימטריות כמו ReLU, שמאפסות קלטים שליליים, He מגדילה את השונות ל-:

$$(12.10) \quad W_{ij} \sim \mathcal{N}\left(0, \frac{2}{n_{\text{in}}}\right)$$

כך נשמר מעבר הגרדיאנטים למרות הדילול שנגרם על ידי ReLU.

סיכום

רשתות עמוקות רגישות לסקאלה ולהתפלגות של אקטיבציות וגרדיאנטים. נורמליזציה - ובמיוחד BatchNorm וLayerNorm - מטפלת בחוסר יציבות על ידי סטנדרטיזציה של ערכי הפיצרים ובכך משפרת את מעבר המידע בין שכבות. במקביל, אתחול נכוון מבטיח שהגרדיאנטים לא יהיו גדולים מדי ולא קטנים מדי בתחילת האימון.

אתחול Xavier אידיאלי לרשתות עם אקטיבציות סימטריות, בעוד אתחול He מותאם למודלים מבוססי ReLU. סכום אלו מזערות עיוות נתונים ומשפרות את יעילות הלמידה.

ביחד, נורמליזציה ואתחול מהווים את עמוד השדרה של דינמיקת אימון יציבה, מהירה ורוביומית בלמידה عمוקה מודרנית. בשיעור הבא נחקרו כיצד להעшир דאטה עבור אימון בצורה יעה (Data Augmentation), כדי לשפר הכללה מבליל לאסוסף דוגמאות נוספת.

13 שיעור

אוגמנטציה

מבוא

מודל עשווי להתאים ביציבות, להתכנס בצורה חלקה ולהציג Training Loss מוקן - ובכל זאת להיכשל בצורה דרמטית במימוש. מדוע זה קורה? בעולם האמתי, DATA הוא רועש, לא עקי ולייטים רוחקות דומה לדוגמאות המטופחות ששימשו אותנו באימון. המפתח להצללה - היכולת של המודל להצליח על קלטים שלא ראה קודם - אינו תמיד חבוי בארכיטקטורה עצמה. לעיתים קרובות הוא טמון בDATA. **Data Augmentation** הוא התהlik של הרחבה והעשרה הדאטאסט על ידי מימוש טרנספורמציות השומרות על נכונות התוויות. הוא מגדיל רובסיטיות על ידי סימולציה של שונות מהעולם האמיתיי וכופה על המודל להפוך לאיננו-ריאנטי לשינויים לא רלוונטיים.

מתי Augmentation עוזר?

Data Augmentation שימושי במיוחד כאשר DATA האימון קטן, לא מאוזן, או חסר גיוון. במקרים כאלה המודל נוטה ל-Overfitting - השגת דיוק גבוהה על DATA האימון אך כישלון בהצללה לדוגמאות חדשות. אוגמנטציה מכניתה וריאציות מבוקרות לDATA ועוזרת למודל להפוך לרוביטי לשינויים בתתפלגות. זה חשוב במיוחד כאשר הדומין של הבדיקה כולל דפוסים או UIותים שאינם מיוצגים היטב בDATA האימון. אסטרטגיית Augmentation

טובה תגרום בדרך כלל לעלייה קלה ב-Training Loss, כי המודל מתמודד עם קלטים קשים יותר ומגוונים יותר, אך טוביל לירידה ניכרת ב-Validation Loss ולשיפור ב-Accuracy על דатаה שהמודל לא ראה. כתוצאה לכך, Augmentation מהויה שיטה פשוטה אוחזקה לשיפור הכללה ואמינות, במיוחד במקרים מסוימים ראייה או כאשר הדאטה המתוג מוגבל.

Generation Augmentation לעומת

חשוב לבחין בין Augmentation ל-Generation. Augmentation אינו מייצר דוגמאות חדשות לחלוטין. במקום זאת, הוא יוצר גרסאות חלופיות של אותה דוגמה תוך שמירה על התוויות הסמנטיות. נניח כי x היא תמונה קלט ו- y התוויות שלה. Generation מייצר גרסה מומרטת (x' כך ש-):

$$(13.1) \quad f(x') \approx f(x), \quad \text{and} \quad \text{label}(x') = y$$

כאשר T היא פונקציית טרנספורמציה שנלקחת מהתפלגות של "шибושים" לגיטימיים ו- f הוא המודל.

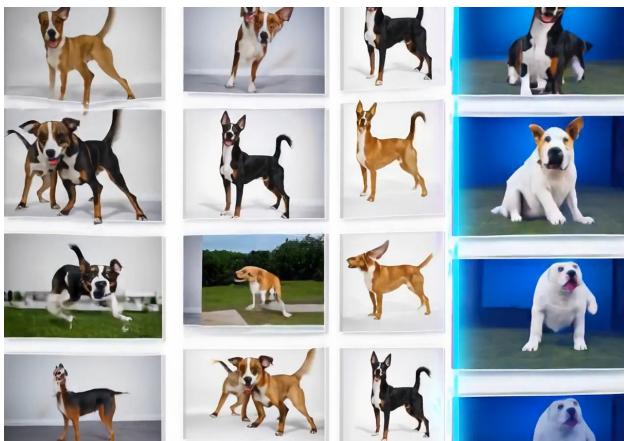
שמירת עקביות התוויות

דרישה בסיסית היא שהדוגמה המומרטת חייבת לשמר על התוויות המקוריות. הפרת כל זה עלולה לגרום לתיג רועש" ולביצועים גרועים. למשל, סיבוב הספרה "6" ב- 180° יוצר משהו שקשה להבחין ביןו לבין "9" - אך סימונו כ-"6" יטעה את המודל. זה כבר לא Augmentation, אלא הרעשה והטעיה סמנטית.

טרנספורמציות נפוצות

בתחום הוייזואלי, נהוגות הטרנספורמציות הבאות:

- $x' = \text{Flip}_H(x)$:**Horizontal Flip** •
- $x' = \text{Rotate}_\theta(x), \theta \in [-15^\circ, 15^\circ]$:**Rotation** •
- :**Crop and Resize** •
- :**Translation and Zoom** •
- .**Mechanisms**.
- :**Brightness and Contrast** •
- :**Color Jitter** •
- :**Additive Noise** •



. איור 13.1: דוגמאות ל-**Augmentation**

זהירות

עודף Augmentation עלול להזיק. אם הדוגמאות המתקבלות כבר אין דומות לקליטים תקפים - למשל עיוות קיצוני או אובדן מבנה אובייקט - הדבר עלול לפגוע בلمידה. בדיקה ויזואלית חיונית.

מיימוש

ספרייה מומלצת ל-Image Augmentation היא **Albumentations**, המשתלבת היטב הhn עם PyTorch והhn עם TensorFlow. ניתן לבצע באופן אקלרי היפוכים, סיבובים, שינוי בהירות וNORMALIZACIA - תוך שמירה על הסמנטיקה של התווית.

Online Augmentation Offline

ישנן שתי אסטרטגיות ליישום Augmentation:

- **Offline Augmentation**: יוצר ושמור דוגמאות חדשות מראש. מגדיל שימוש בזיכרון אך מספק וריאציה קבועה.
- **Online Augmentation**: מיישם טרנספורמציות אקראיות בזמן אמת במהלך האימון. מספק סטטיסטיות וגיון אינסופי לאורך אפוקים.

הרבית המערכות המודרניות מעדייפות Online Augmentation בשל הגמישות ויעילות הזיכרון.

Best Practices

כדי למקסם את האפקטיביות של Augmentation:

- .1 **Validate visually**: תמיד בדקו באצ' של תמונות מומורות.
- .2 **Prioritize plausibility**: הימנעו מטרנספורמציות שמספרות זהות מחלוקת.
- .3 **Prefer small, composable transformations**: עדיף לשדרש פעולות פשוטות מאשר לבצע עיוזת אגרסיבי יחיד.
- .4 **Adapt per domain**: Augmentation למספרות אינו זהה לאלה לתמונות רפואיות או לתמורי תנועה.

סיכום

Augmentation היא אחת הטכניקות החזקות והנגישות ביותר לשיפור רובוטיות של מודלים. היא מדמה שנות, כופה אינטראקטיות ומחיתה תלות ברגולרייזציה דנית. בנגדן לשינויים בארכיטקטורה או לפונקציות מחיר חדשות, אוגמנטציה משנה את הדטה עצמה - וועזרת למודל לפגוש את העולם כפי שהוא באמת: רועש, מגוון ובלתי צפוי.

אוגמנטציה מותוכנתה היטב אינה רק מעשירה את DATASET האימון - אלא מרחיבה את נקודת המבט של המודל. היא מכינה את הרשות לא לשנן, אלא להבין.

שיעור 14

אופטימיזציה מתקדמת

למידה עמוקה מודרנית לא הייתה אפשרית ללא אלגוריתמי אופטימיזציה יעילים ויציבים. בעוד שרשתות נוירוניות עשוות להיראות טובות בתחזיותיהן, בלביתן הן נשענות על עדכנים חזריים זהירים למיליאוני פרמטרים. בשיעור זה נבחן כיצד מתבצעים עדכנים אלו - מהיקרונו הבסיסי של Gradient Descent ועד האופטימיזרים המתקדמים ביותר כיום: AdamW, Lion, Adafactor ועוד. אופטימיזרים אלו קובעים כיצד המודל לומד, באיזו מהירות הוא מתכנס והאם הוא מצליח להקליל.

מ-*Backpropagation* לאופטימיזציה

כל צעד אימון מתחילה ב-*Backpropagation*, שמחساب את הגרדיאנטים של פונקציית המחיר ביחס לכל משקל בראשת. גראדיאנטים אלו מצביעים על הכיוון שבו המחיר גדל ביותר. כדי לשפר את המודל, אנו רוצחים ללכת בכיוון הפוך (מצערו המחיר). האופטימיזר הוא האלגוריתם שהופך גראדיאנטים לעדכוני פרמטרים ממשמעותיים.

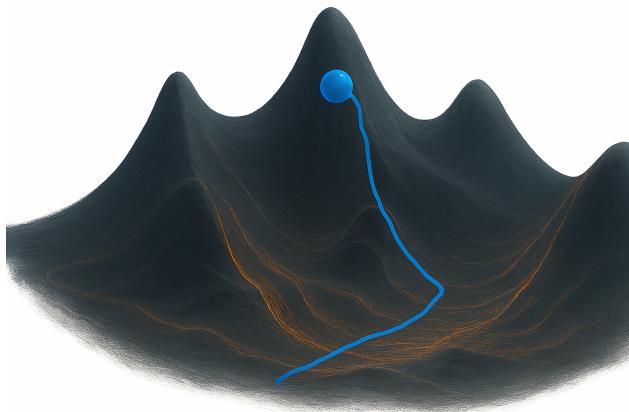
חזרה לבסיס: Gradient Descent

כפי שראינו קודם, כלל העדכון הקליני פשוט ויסודי:

$$(14.1) \quad \theta_{k+1} = \theta_k - \eta \cdot \nabla_{\theta} L(\theta_k)$$

כאשר θ_k מייצג את וקטור הפרמטרים בצעד k , η הוא קצב הלמידה ו- $\nabla_{\theta} L(\theta_k)$ הוא הגרדיאנט של פונקציית המחיר L בנקודה זו. בambilים פשוטות - המודל "גולש" על פני משטח המחיר, מונחה על ידי השיפוע המקומי. העונה: זו אותה משווהה שראינו קודם, אך כאן המשקל מסומן ב- θ במקומות w .

בפועל, משטחי אימון אינם חלקים; הם רועשים, עתירי ממדים ומלאים ב"עמקים" ו"ערוצים". Gradient Descent פשוט מתנסה להתמודד עם מרכיבות זו - וכן יש צורך באופטימיזרים מתקדמים יותר.



איור 14.1: אלגוריתם Gradient Descent של: כדור מתגלגל על פני משטח מחיר מורכב, מונחה על ידי השיפוע התולול ביותר בכל נקודה. המסלול מייצג את תהליך האופטימיזציה כאשר המודל מתאים את פרמטריו כדי למזער את השגיאה. פסגות ועמקים מדגימים את האתגר ביציאה מנקודות מינימום מקומיות לעבר מינימום גלובלי.

Adam: Adaptive Moments

ידי שילוב שני מנגנונים: Adam - Adaptive Moment Estimation וסמיילינג אדפטיבי. הוא שומר ממוצע

מעריצי דועך של הגרדיאנטים (מומנט ראשון) ושל הריבוע שלהם (מומנט שני):

$$(14.2) \quad m_k = \beta_1 \cdot m_{k-1} + (1 - \beta_1) \cdot g_k$$

$$(14.3) \quad v_k = \beta_2 \cdot v_{k-1} + (1 - \beta_2) \cdot g_k^2$$

$$(14.4) \quad \theta_{k+1} = \theta_k - \eta \cdot \frac{m_k}{\sqrt{v_k} + \epsilon}$$

כאשר g_k הוא הגרדיאנט בצעד k ו- m_k , v_k הם ממוצעים נועים של הגרדיאנט ושל ריבועו. האיבר ϵ (בדרך כלל 10^{-8}) מבטיח יציבות נומרית. ברירות מחדל מקובולות: $\beta_2 = 0.999$, $\beta_1 = 0.9$.

Adam מתאים את גודל הצעד לכל פרמטר בנפרד, מה שהופך אותו לרובטיבי ויעיל במידה רחבה של משימות. הוא עדין נחשב לאופטימייזר ברירת המחדל אצל רבים.

AdamW: Decoupled Weight Decay

AdamW הוצג כדי לתקן בעיה עדינה במימוש הרגולרייזציה של Adam. ב-AdamW מפריד ביןיהם ומימוש Weight Decay (L2 Regularization) ב�וראה לא רצואה עם קצב הצעד Weight Decay ב�וראה מפורשת: האדפטיביים.

$$(14.5) \quad \theta_{k+1} = \theta_k - \eta \cdot \left(\frac{m_k}{\sqrt{v_k} + \epsilon} + \lambda \cdot \theta_k \right)$$

כאן, λ הוא מקדם Weight Decay (למשל 0.01). ניסוח זה משפר הכללה יציבות באימון, במיוחד במקרים גדולים מבוססי Transformers. כתוצאה לכך, AdamW הוא האופטימייזר המעודף בארכיטקטורות מתקדמות רבות.

Lion: Lightweight Momentum with Signs

Lion, שהוצג בשנת 2023 על ידי חוקרים בגוגל, מציע טויסט אלגנטיבי: במקום

להשתמש בעוצמת הגרדיאנטס או במומנטים שניים, הוא מעדכן משקלות על בסיס סימן המומנטום בלבד:

$$(14.6) \quad m_k = \beta_1 \cdot m_{k-1} + (1 - \beta_1) \cdot g_k$$

$$(14.7) \quad \theta_{k+1} = \theta_k - \eta \cdot \text{sign}(m_k)$$

גישה זו מקטינה שימוש בזכרון ועלות חישובית, תוך שמירה על הכוון. על במיוחד במודלים ויואליים כמו ViTs והוכח כיעיל W-Lion על אימון מסוימים, במיוחד כאשר זיכרון ה-UPS מוגבל. היפר-פרמטרים מקובלים: $\eta = 3 \times 10^{-4}$, $\beta_1 = 0.9$, $\beta_2 = 0.99$.

יעילות זיכרון בקנה מידת גדול: Adafactor

פותח עבור מודלים עצומים, שבהם אחסון מומנטום שניים מלאים הופך לצוואר בזיכרון. במקום לשמר מטריצות מלאות, Adafactor מפרק אותן לוקטורים של שורות ועמודות - מצמצם את דרישת הזיכרון מ- $O(n^2)$ ל- $O(n)$. כך ניתן לאמן מודלי שפה ענקיים (כמו T5) גם על חומרה מוגבלת. Adafactor תומך גם במנגנון קצב למידה אדפטיבי לפי סקליל הפרמטרים וניתן לשלבו עם שיטות זמן צמוי כמו Cosine Decay.

קצב למידה וזמן אופטימיזציה

האופטימיזר הוא רק חלק מההמונה; זמן של קצב למידה חשוב לא פחות. קצב למידה קבוע כמעט אף פעם לא מתאים לכל שלבי האימון. אסטרטגיות נפוצות כוללות:

Warmup + Decay: התחלת בקצב למידה נמוך, עלייה הדרגתית ואז דעיכה.

Cosine Annealing: הפחיתה הדרגתית של הקצב לפי עקומת קוסינוס, לרוב יחד עם Warmup.

AdamW ו-Lion מציגים ביצועים טובים עם תזמון של קצב הלמידה לפי פונקציית קוסינוס, בעוד ש-Adafactor משולב לעיתים עם תזומנים אדפטיביים לפי נורמות פרמטרים. בחירת תזמון נכון יכולה להשפיע משמעותית על מהירות ההתקנסות ועל הביצועים.

שיקולי יציבות והכללה

אופטימיזרים אדפטיביים כמו Adam ו-Lion מתכנסים מהר, אך עלולים להיקלע ל-Minima Sharp - נקודות במרחב הפרמטרים שמצוירות מחיר אימון אך מכליות גרען. SGD עם מומנטום, על אף אייטוינו, נוטה למצוא נקודות מינימום נמוכות יותר - מה שМОוביל לדיווק טוב יותר. בתנאים רועשים במיוחד, שיטות אדפטיביות (ובעיקר AdamW) נוטות להיות יציבות יותר. עם זאת, מעבר לאופטימיזרים פשוטים כמו SGD בשלב Fine-Tuning עשוי לשפר רובייטיות.

טבלת השוואת

אופטימיזר	זיכרון	קיוב למידה אדפטיבי	מומנטום	שימוש אופfineyi
SGD	נמוך	לא	כן	משימות רובייטיות
AdamW	בינוני	כן	כן	רוב המודלים
Lion	נמוך	לא	סימן	ViTs, מעט זיכרון
Adafactor	נמוך מאוד	כן	לא	מודלים עצומים

טבלה 14.1: סיכום אופטימיזרים מודרניים

טיפים

עבור דאטסהטים רועשים או גרדיאנטים לא יציבים, AdamW הוא בחירה בטוחה ויעילה. עבור סביבות GPU מוגבלות - במיוחד באימון ViTs - עשוי לספק עילوت גבוהה יותר. במקרים גדולים מאוד (למשל מעל מיליארד פרמטרים), Adafactor הופך חיוני בזכות עילותו.

סיכום

האופטימיזרים המודרניים הם הרבה יותר מ-Gradient Descent פשוט. הם משלבים רעיונות של מומנטום, קצבי למידה אדפטיביים, עילות זיכרון ורגולרייזציה - כל אחד מהם מציע יתרונות בתרחיש שונה.

AdamW אומץ באופן נרחב בזכות אייזון בין יציבות לביצועים; Lion מציע פשטות ומהירות במודלים ויוזאלים תחת מגבלות משאבים; Adafactor מאפשר אימון בקנה מידה עצום; ו-DGD, על אף פשטותו, נותריעיל בעיקר בשלבי Fine-Tuning.

שיעור 15

להבין החלטות של מודלים Explainability

מבוא

כל מודלים של במידה עמוקה משמשים יותר ויותר בקבלת החלטות בעלות סיכון גבוה, היכולת להסביר את פעולותם הופכת להכרחית. בתחוםים כמו בריאות, פיננסים, מערכות אוטונומיות וביתחון - כבר לא מספיק שמודל יהיה מדויק; עליו גם להיות מובן.

Explainability מתייחסת למגוון שיטות המספקות תובנה מדויקת יוצר פלט מסוים. לשיטות אלו מטרות רבות: הן מבירות אמון משתמשים, עוזרות למפתחים לנפות התנהגויות בלתי צפויות, מבטיחות עמידה ברגולציה וחושפות סוגיות אטיות או של הוגנות. בסופו של דבר, מודל שלא ניתן להסביר - הוא מודל שיהיה קשה לסמוך עליו.

מדוע קשה להסביר רשותות עמוקות

מודלים ליניאריים ועצי החלטה מספקים שקיפות מובנית: ניתן לבדוק ישירות מקדמים או כללי החלטה. לעומת זאת, רשותות נוירוניים עמוקות פועלות באמצעות שכבות רבות של טרנספורמציות לא-ליניאריות, לרוב עם מיליון פרמטרים. הייצוגים הפנימיים שלחן מבזרים, שאזורים ובדרכן כלל בלתי ניתנים לפרשנות ללא ניתוח Post Hoc.

מורכבות זו מתקשה לקבוע איזה מידע השפיע על ההחלטה מסוימת, או

לזהות מתי המודל מסתמך על קורלציות מקריות או משתנים מעורבים.

מקרה לדוגמה: גילוי הונאות בبنקרים

נניח בנק המפעיל מודל למידה عمוקה לגילוי עסקאות של הונאת אשראי. במהלך הבדיקות המודל מציג דיקוגרף גבוה, אך צוות הרגולציה דורש שכל עסקה חשודה תלווה בהסביר. כאנליסטים חוקרם מספר False Positives,

הם מגלים תבנית מט裏יה: הונאות רבות מרוכזות סביב עיר מסוימת. באמצעות הסברי SHAP, הם מגלים שהתכונה "מקום עסקה" תורמת רבות להחלטה - אך רק כאשר היא משולבת עם קטגוריות מסוימות. הדבר מוביל לביקורת עמוקה יותר, החושפת שהמודל למד קשר פעילות "חשודה" לשילוב של התנהוגיות לגיטימיות לחלוין.

במקביל, אנליסטים מפעילים Grad-CAM על תת-מודול מבוסס CNN שמנתח קובלות סרווקות. Feature Maps מגנות שהמודל מתמקד לעיתים בטקסט מודגש - ולא בפרטיו העסקי - בשל הטוויות בדעתה הסט האימון. שתי שיטות ההסבר הללו, שיישמו ברמות שונות במערכת, סייעו לצוותאמן מחדש את המודל, לשפר את תהליכי הדאטה ולעמוד בדרישות הרגולטוריות לשקיפות.

CNNs: הסבר חזותי ב-

Grad-CAM

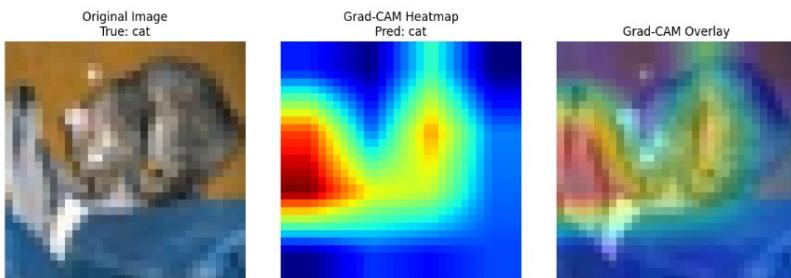
ברשותם קונבולוציה, כלי סטנדרטי להסביר פועלות המודל הוא Grad-CAM (Gradient-weighted Class Activation Mapping). מדגישה את האזוריים המרחביים בתמונה שהשפיעו ביותר על החלטת המודל. נסמן ב- y^c את ציון המודל (לפניהם Softmax) עבור מחלקה c ונסמן $A^k \in \mathbb{R}^{H \times W}$ כמפת הפיצרים $-k$ מהשכבה הקונבולוציונית שנבחרה. החישובות α_k^c מחושبات כך:

$$(15.1) \quad \alpha_k^c = \frac{1}{Z} \sum_{i=1}^H \sum_{j=1}^W \frac{\partial y^c}{\partial A_{i,j}^k}$$

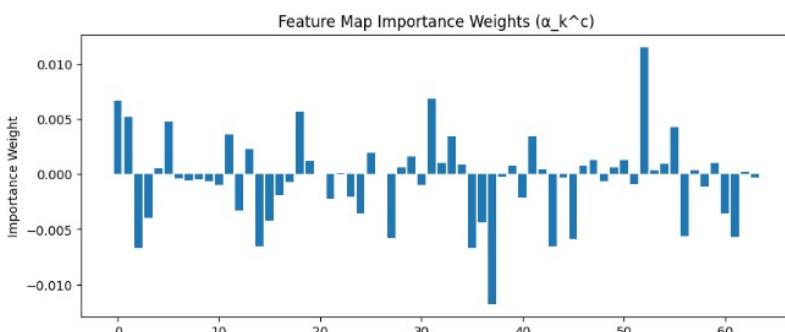
כאשר $Z = H \cdot W$ הוא מספר האלמנטים המרחביים במפת הפיצ'רים. משקלים אלו מייצגים עד כמה כל מפת פיצ'רים משפיעה על ציון המחלקה. מפת ה- h-CAM עצמה מחושבת כך:

$$(15.2) \quad L_c^{\text{Grad-CAM}} = \text{ReLU} \left(\sum_k \alpha_k^c A^k \right)$$

הפעלת ReLU מבטיחה שרק אזורים שטורמים ערכיהם חיוביים. מפת החום מוגדלת חוזרת לגודל הקלט לצורך ויזואלייזציה.



איור 15.1: דוגמה - המחתת ניבוי "חתול" ב-10-CIFAR. התמונה מציגה הסבר עבור CNN שאומן על CIFAR-10. המודל סיוג נכון התמונה כ-*cat* ו-*CAM* מדגיש את האזורים הרלוונטיים ביותר שהובילו לניבוי.



איור 15.2: החשיבות α_k^c עבור כל מפת פיצ'רים. כל מפת פיצ'רים לומדת תבנית ויזואלית מסוימת של ה- h-CNN (למשל טקסטורות, צורות, חלקיים). משקלים חיוביים תומכים בסיווג "חתול"; שליליים סוננו על ידי ReLU.

חשיבות גלובלית: ניתוח מבוסס Permutation

במודלים טבלאיים, הפרשנות מתמקדת לעיתים בזיהוי אילו פיצ'רים חשובים באופן גלובלי. בעוד שמודלים מבוססי עצים (כמו XGBoost) מספקים מדריכים פנימיים, גישה אינטואיטיבית ואגנושטיית-למודל היא **Permutation Importance**. נניח f הוא מודל מאומן ו- \mathcal{D} DATA וידציה. עבור כל פיצ'ר x_i , מחליפים את ערכיו אקראית ומודדים את שינוי הביצועים. פורמלית:

$$(15.3) \quad \Delta_i = M_0 - M_i$$

כאשר M_0 הוא המדריך המקורי (למשל דיקון) ו- M_i המדריך לאחר שינוי x_i . ערך גדול של Δ_i מצביע על כך שהפיצ'ר חיוני לביצועים.

הסבירים מקומיים מתמטיים SHAP:

SHAP (SHapley Additive exPlanations) היא שיטה לכימות תרומת פיצ'רים לפלט המודל. היא מבוססת על ערכי Shapley מתוך המשחקים ומקיימת ארבעה אקסימוטות: יעילות, סימטריה, דמה ואדיטיביות - המבטיחות ייחוס הוגן ועקבני.

נסמן F כקובצת הפיצ'רים המלאה ו- (S) את ניבוי המודל באמצעות תת-קובוצה S . ערך ה-SHAP עבור פיצ'ר i מוגדר כך:

$$(15.4) \quad \phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|!(|F| - |S| - 1)!}{|F|!} [f(S \cup \{i\}) - f(S)]$$

כלומר:

- ϕ_i הוא ציון הייחוס של פיצ'ר i .
- $f(S)$ הוא פלט המודל הצפוי עם תת-קובצת הפיצ'רים S .
- הפקטור הקומבינטורי משקל את כל הצירופים האפשריים של הפיצ'רים.

בפועל, חישוב מדויק של משווהה זו יקר מאד עבור פיצרים רבים. לכן משתמשים בקירובים כגון Kernel SHAP (אגנוטטי למודל) או Deep SHAP (לרשתות נוירוניות). אלו משתמשים בהתפלגות ידועות או דוגמאות יייחוס. למרות העלות, SHAP מספק הסברים יציבים אינטואיטיביים ברמת הדוגמה הבודדת והפץ לכלי מרכזי בתעשיות מופוקחות.

הסבר טקסט באמצעות Attention Weights

במודלים מבוססי Transformers כגון BERT או GPT, פרשנות מוטמעת כבר בארכיטקטורה באמצעות מנגנון $\text{-} \text{Self-Attention}$. כל טוקן מחושב כצירוף משוקלל של כל שאר הטוקנים, ליצירת ייצוג הקשור. נזכיר את נוסחת $\text{-} \text{Attention}$ המשיעורים הקודמים:

$$(15.5) \quad \text{Attention}(i, j) = \frac{\exp(q_i^\top k_j / \sqrt{d_k})}{\sum_{j'=1}^T \exp(q_i^\top k_{j'} / \sqrt{d_k})}$$

כאשר:

- $q_i, k_j \in \mathbb{R}^{d_k}$ הם וקטורי $\text{-} \text{Query}$ וה- $\text{-} \text{Key}$ של טוקנים i ו- j .
- d_k הוא ממד מרחב $\text{-} \text{Key}/\text{-} \text{Query}$.
- T הוא מספר הטוקנים בדף.

משקלים $\text{-} \text{Attention}$ יוצרים מטריצה $A \in \mathbb{R}^{T \times T}$ שנייה להציג כדי לגלוות אילו מילים משפיעות זו על זו בכל שכבה וראש. למרות $\text{-} \text{Attention}$ אינו מספק הסבר ישיר לניבוי, הוא חשוב את אופן פעלת המודול בהבנת הקשר, במיוחד במקרים של תרגום או ניתוח תחבירי.

מגבליות ואתגרים

למרות תועלתן, לשיטות פרשנות מגבלות שונות: Grad-CAM מוגבל בראזולציה ועלול לייצר מफות רועשות בשכבות عمוקות. SHAP, על אף היסודות המתמטיים,

רגיש לקורלציה בין פיצ'רים ויקר חשיבות בימיד גובה. Attention Weights משקפים אסוציאציה ולא סיבתיות ויכולים להשתנות בין ראשיים וশכבות. הסברים גם עלולים להטעות: מודל עשוי ליחס ערכי SHAP גבוהים לפיצ'רים שהם רק "פרוקסים" למשתנים אחרים - עשויים להציג טוקנים תחביריים ללא רלוונטיות סמנטית. לכן, צריך להיחס עדה דיאגנositית - לא מקור אמת מוחלט. שילוב עם ידע תחומי, בדיקות רובוטיות והערכת אונושית - חיוניים.

סיכום

הופך במידה עמוקה שchorה למערכת הניננת לביקורת Explainability ולפועלה. Grad-CAM מציג אזורים חזותיים משמעתיים ב-CNNs. SHAP מספק "יחס מתמטי" מבוססת תורת המשחקים לניבויים אינדיבידואליים. שיטות Permutation מספקות תובנות גLOBליות פשוטות אך אפקטיביות. Attention Weights פותחים חלון להבנת אינטראקציות ברמת טוκן במודלים מבוססי Transformers לכל טכניקה יתרונות ומוגבלות. הבחירה ביניהן תלויות בסוג הדאטה, המודל ומטרות הפרשנות.

שיעור 16

PyTorch מול TensorFlow

מבוא

לאחר שלמדנו כיצד לתוכנן, לאמן ולמשרשות נוירונים, עליה לעתים שאלת פרקטית:izia פוריימוורק להשתמש? שתי הספריות הפופולריות ביותר כיום הן TensorFlow ו-PyTorch. שתיהן בשימוש נרחב, מתחזקות היבט ותומכות הן בסביבות מחקר והן בסביבות פרודקشن. הן חולקות הרבה מהפונקציונליות - אך נבדלות בפילוסופיה, בתחביר ובדפוסי שימוש אופייניים.

בשיעור זה נשווה בין שתי הספריות מבחינה מבנה, סגנון ואקוסיסטים. בנוסף נציג דוגמת סיווג עם Fashion MNIST כדי להמחיש את ההבדלים בקורס.

רקע היסטורי ואקוסיסטים

הוזג על ידי גугл בשנת 2015 כפלטפורמה סקלאbilית מוכוונת TensorFlow. הוא הפך לברירת מחדל עבור יישומים ארגוניים רבים, בזכות כלים כמו TensorFlow Lite, TensorFlow Serving ו-TensorFlow.js. PyTorch הושק על ידי פיזיובוק בשנת 2016 ואכח ב מהירות לפופולריות בקהילה המדעית. התחביר הפיתוני הטבעי שלו, מודל ההערכת הדינמי והגמישות הפכו אותו לבחירה מועדפת לניסויים ולפרוטוטייפים מהירים. עם הזמן, שתי הספריות התפתחו. TensorFlow 2 נמצא ברירת מחדל

של TorchServe ו-PyTorch כוללת כו"ם כלים למימוש, כמו Eager Execution ו-ONNX (Open Neural Network Exchange) ו-TorchScript. שתיהן תומכות ב-ML מודרניים. ומשתלבות היטב עם כלי ML מודרניים.

מודלי הרצה: סטטי מול דינמי

ההיסטוריה, TensorFlow דרשה מהמשתמש להגדיר גרפ' חישוב סטטי מראש, שהורץ בתוך Session. מודל זה הקל על אופטימיזציה אך צמצם גמישות. לעומת זאת, PyTorch השתמשה מלכתחילה בחישוב דינמי: הקוד מושך שורה אחר שורה כפי שנכתב בפייתון. זה מקל על דיבוג ובניאת מודלים עם מנגנון בקרה מורכב.

ב-2 TensorFlow, בירית המclid היא Eager Execution. יחד עם זאת, ניתן להשתמש ב- @tf.function כדי להמיר פונקציות פיותון לייצוגי גרף - שילוב של גמישות וביצועים.

מקרה לדוגמה: סיוג Fashion MNIST

כעת נדגים שימוש עם שתי הספריות על משימות סיוג זהה: Fashion MNIST. הדאטהסט כולל תמונות גווני אפור (28×28 פיקסלים) של פריטי לבוש, מחולקות ל-10 מחלקות. המטרה: למשר רשת נוירונית פשוטה לשיווג.

מימוש ב-PyTorch

```
import torch
import torch.nn as nn
from torchvision import datasets, transforms
from torch.utils.data import DataLoader

train_data = datasets.FashionMNIST('.', train=True, download=True, transform=
    transforms.ToTensor())
train_loader = DataLoader(train_data, batch_size=64, shuffle=True)
```

```

class Net(nn.Module):
    def __init__(self):
        super().__init__()
        self.net = nn.Sequential(
            nn.Flatten(),
            nn.Linear(784, 128),
            nn.ReLU(),
            nn.Linear(128, 10)
        )

    def forward(self, x):
        return self.net(x)

model = Net()
optimizer = torch.optim.Adam(model.parameters(), lr=0.001)
loss_fn = nn.CrossEntropyLoss()

for x, y in train_loader:
    pred = model(x)
    loss = loss_fn(pred, y)
    loss.backward()
    optimizer.step()
    optimizer.zero_grad()

```

שימוש ב-TensorFlow

באמצעות TensorFlow, Keras מציעה מסך דקלרטיבי יותר. היא מפשטת את לולאת האימון ואת הלוגיקה הפנימית, מה שהופך אותה לאידיאלית לניסויים מהירים או למפתחים חדשים.

תחביר, מודולריות ורמת הפשטה

הקוד מעלה ממחיש את ההבדל הסגוני המרכז:

- PyTorch מדגישה שקייפות ומודולריות. המשתמש מגדיר במדויק את המודל, את חישוב המידע בראש וראשיתelogיקת האופטימיזציה.
- TensorFlow עם Keras מפשטת פרטיים אלה, מאפשרת פיתוח מהיר יותר, אך דורשת קונפיגורציה עמוקה יותר במשימות מתקדמות.
- במהדרת PyTorch מרגיש כמו פיתון טבעי, בעוד שקוד TensorFlow - במיוחד kod framework מודלים פונקציונליים או במצב גוף - מרגיש יותר כמו עם מוסכמות משלו.

אקויסיטטס ופרויקטים

בבחירה framework יש להתחשב לא רק בפיתוח אלא גם בIMPLEMENTATION:

- TensorFlow תומך ב:
- TensorFlow Serving למימוש ב-REST ו- gRPC
- TensorFlow Lite להרצה במובייל/אמביד
- TensorFlow.js להרצה בדפדפן
- PyTorch תומך ב:
- TorchScript לייצוא מודלים כגרפים סריאליים
- TorchServe לפריסת מודלים
- ONNX, המאפשר שימוש במסגרות אחרות
- Mixed-Precision, האצה ב-GPU/TPU ואינטגרציה
- .MLflow Weights & Biases ,Hugging Face עם

מתי להשתמש בכל אחת

הבחירה תלויה בהקשר. PyTorch מצוינת לניסויים מהירים, מחקר והוראה. היא מספקת שליטה מדויקת ודיבוג ברור. TensorFlow מצוינה בIMPLEMENTATION קצח-לקצה, מערכות פרויקט-טוווח ופריסות מובייל/ווב.onznx ביכולת תמיכת ONNX, ניתן להעיבר מודלים ביניהם. נפוץ יותר ויותר לבצע פרוטוטיפ PyTorch ולהמיר ל-TensorFlow פרויקט - או להפוך.

סיכום

גם TensorFlow וגם PyTorch הן ספריות מודרניות עצמאיות. הן חולקות תכונות רבות ודוגמאות בתוכן. הבדלים מרכזיים נמצאים בתחריב, בرمת ההפשטה ובכלי האקוסיטם. בחרו ב-PyTorch לשליתה, גמישות ומידול דינמי. בחרו ב-(TensorFlow (Keras

שיעור 17

עבודה אפקטיבית עם Google Colab

מבוא

שיעור זה מציג את תהליך העבודה ב-*Colab*, כולל הדרות חומרה, ניהול קבצים, טיפים להתקנות, אימון מודלים ופרקטיקות מומלצות להימנעות מניתוקים ולמكسום פרודוקטיביות. שלו.

שיעור זה מציג את תהליך העבודה ב-*Colab*, כולל הדרות חומרה, ניהול קבצים, טיפים להתקנות, אימון מודלים ופרקטיקות מומלצות להימנעות מניתוקים ולמקסום פרודוקטיביות.

תחילת עבודה בסביבת Colab

כאשר פותחים מחברת חדשה ב-*Colab*, עובדים על Jupyter Notebook מבוסס ענן הרץ על התשתיות של גугл. כל תא במחברת יכול להכיל קוד, הסברים ב-*Markdown*, או תוצאות. הקוד מורץ במכונה וירטואלית מבודדת ומשאים כמו זיכרון, דיסק ושימוש ב-GPU מוצגים בפינה הימנית העליונה של המשק. ניתן לעבור בין בלוקי קוד לטקסט באמצעות הכפתורים "Code" ו- "Text" בסרגל העליון.

```

File Edit View Insert Runtime Tools Help
Commands + Code + Text > Run all
import tensorflow as tf
import numpy as np
import matplotlib.pyplot as plt
from tensorflow.keras import layers, models
import cv2

# Load and prepare data
print("Loading CIFAR-10 dataset...")
(x_train, y_train), (x_test, y_test) = tf.keras.datasets.cifar10.load_data()
x_train, x_test = x_train / 255.0, x_test / 255.0

class_names = ['airplane', 'automobile', 'bird', 'cat', 'deer',
               'dog', 'frog', 'horse', 'ship', 'truck']

# Create and train model using Functional API
print("Creating simple CNN...")
inputs = layers.Input((32, 32, 3))
x = layers.Conv2D(32, (3, 3), activation='relu')(inputs)
x = layers.MaxPooling2D((2, 2))(x)
x = layers.Conv2D(64, (3, 3), activation='relu')(x)
x = layers.MaxPooling2D((2, 2))(x)
conv_output = layers.Conv2D(64, (3, 3), activation='relu', name='target_conv')(x)

# Add a linear layer at the end
x = layers.Flatten()(conv_output)
x = layers.Dense(64, activation='relu')(x)
outputs = layers.Dense(10, activation='softmax')(x)

model = tf.keras.Model(inputs=inputs, outputs=outputs)
model.compile(optimizer='adam', loss='sparse_categorical_crossentropy', metrics=['accuracy'])

model.fit(x_train, y_train, epochs=10, validation_data=(x_test, y_test))

model.evaluate(x_test, y_test)

```

. איור 17.1: צילום מסך של Colab Google

הפעלת מאיצי GPU או TPU

מאפשר לבקש גישה ל-GPU או TPU, מה שמאיץ משמעותית את האימון.

1. גשו ל-`Runtime` < Runtime

2. בחרו GPU או TPU

3. לחצו Save

לאחר הפעלה, ניתן לאמתת את הגישה ל-GPU באמצעות:

```
import torch
print(torch.cuda.is_available())
```

או על ידי בדיקת התחkon הפעיל:

```
!nvidia-smi
```

פלט נפוץ כולל התקנים כמו Tesla T4 או V100.

התקנת חבילות Python

תומך בהתקנת חבילות ישירות באמצעות `pip` בתאי המחברת. לדוגמה:

```
!pip install transformers --upgrade --quiet
```

הדגל quiet – מצמצם לוגים ארוכים של ההתקנה.

חיבור ל- Google Drive

כדי לשמר מודלים, פלטים או נתונים בין שנים, יש להגדיר את Google Drive:

```
from google.colab import drive  
drive.mount('/content/drive')
```

הקבצים יהיו זמינים תחת /content/drive/MyDrive/. כך מבטיחים שהתוכ�ות לא יאבדו כאשר הסשן מתאפשר.

העלאת קבצים ידנית

- ניתן להעלות קבצים מהמחשב המקומי על ידי:
- גירירה אל חלון ה-”Files Pane”
 - שימוש בפקודה:

```
from google.colab import files  
files.upload()
```

קבצים אלו נשמרים זמנית במערכת הקבצים המקומיות של המחברת.

שיתוף ועבודה משותפת

כל מחברת Colab ניתנת לשיתוף כמו מסמך Google Doc. לחזו על כפתור “Share” בפינה הימנית العليا ובחירה הרשאות. ניתן גם לפתח מחברות GitHub מ-

אימון רשת נוירוניים עם GPU

נאמן רשת פשוטה עם PyTorch באמצעות Fashion MNIST ותמיכה ב-GPU:

```
import torch
import torchvision
import torchvision.transforms as transforms
from torch import nn, optim
from tqdm import tqdm

device = 'cuda' if torch.cuda.is_available() else 'cpu'

train_data = torchvision.datasets.FashionMNIST(root='.', train=True, download=
    True,
    transform=transforms.ToTensor())
train_loader = torch.utils.data.DataLoader(train_data, batch_size=64, shuffle=True
    )

model = nn.Sequential(
    nn.Flatten(),
    nn.Linear(784, 128),
    nn.ReLU(),
    nn.Linear(128, 10)
).to(device)

loss_fn = nn.CrossEntropyLoss()
optimizer = optim.Adam(model.parameters())

for epoch in range(3):
    total_loss = 0
    for x, y in tqdm(train_loader):
        x, y = x.to(device), y.to(device)
        pred = model(x)
        loss = loss_fn(pred, y)
```

```

optimizer.zero_grad()
loss.backward()
optimizer.step()
total_loss += loss.item()

print(f"Epoch {epoch}1+, Loss: {total_loss:.f}")

```

ניטור עם TensorBoard

כדי לנטר את האימון באופן ויזואלי, Colab תומך ב-.TensorBoard. תחילה יש לטענו את ההרחבה:

```
%load_ext tensorboard
%tensorboard --logdir runs
```

בתוך לולאת האימון, יש להגדיר את המדדים באמצעות:

```
from torch.utils.tensorboard import SummaryWriter
writer = SummaryWriter()

# האימון במחלך
writer.add_scalar('Loss/train', loss.item(), step)
```

טיפים מתקדמים

- שבירת Model Checkpoints :

```
torch.save(model.state_dict(), 'model.pth')
```

- בדיקת שימוש בדיסק:

```
!du -sh /content
```

- הציגת קבצים:

```
!ls -lh
```

- מדידת זמן ריצה של תא:

```
%time
```

- בדיקת שימוש בזיכרון:

```
!free -h
```

סיכום

היא פלטפורמה נוחה, נגישה עם שימוש ב-UPG למידה عمוקה. Google Colab עם הגדרת עבודה נכון, הוא הופך ליותר מלווה ניסויים - היא הופכת למחשבה מלאה בענן. למדו את הקיצורים הקיימים בסביבה לעבודה עיליה, נצלו בחכמה את ה-Drive, נטרו את האימון באופן ויזואלי ועבדו במבנה מודולרי.

שיעור 18

Mixed Precision Training

מבוא

כעת ראיינו כיצד בונים, מאמנים, מאפטמים ומעריכים רשותות עמוקות. אך יש שכבה אחורונה שמחברת את כל אלה - ממד נסתר שבו זמן, זיכרון ודיקום מתלכדים.

Mixed Precision Training אינו אופטימייזר חדש, לא טרייק רגולרייזציה ולא פונקציית מחיר חדשה. זה שינוי יסודי באופן שבו אנו מייצגים ומבצעים חישובים בתוך הרשותות - ואם משתמשים בו נכון, הוא יכול להוביל לאימון מהיר יותר, שימוש מופחת בזכרון ולפעמים אפילו להקללה טובה יותר.

מה המשמעות של Precision בפועל

רוב הרשותות הנוירוגניות כיומ מאומנות באמצעות מספרים בפורמט נקודה צפה של 32 בית (FP32). זה הפורמט המקביל לאחסון משקלות, חישוב אקטיבציות ומעקב אחרי גרדיאנטים. הוא מציע יציבות נומרית גבוהה, אך במחיר: כל מספר צריך 32 בית וכל פעולה משתמש בנתיב חישוב יקר יחסית.

GPUs מודרניים, במיוחד מסדרות Turing ו-Ampere של NVIDIA, תומכים באלטרנטיבה: מספרים בפורמט 16 בית נקודה צפה (FP16). אלו קטנים יותר, מהירים יותר ויעילים יותר בזכרון - אך מקרים מסוימת של דיקום

נומי. ההבדל אינו רק באחסון. FP32 משתמש ב-8 בית לאקספוננט ו-23 למנטייה, בעוד FP16 משתמש ב-5 בית לאקספוננט ו-10 למנטייה. טווח מצומצם זה עלול להוביל לאי-יציבות אם לא מטופל כראוי. למורות זאת, שימוש ב-FP16 יכול להכפיל את רוחב הפס של האזכור ולהפחית משמעותית את זמן האימון. האתגר הוא להחליט היכן ניתן להשתמש בו באופן בטוח - והיכן לא.

הגישה של Mixed Precision

במקום לבחור בין FP32 ל-FP16, צינורות אימון מודרניים משלבים את שניהם. זהו הרעיון של Mixed Precision. כברירת מחדל, חישובים שאינס רגיסטים לדיקוק נמוך מבוצעים ב-FP16. חישובים רגיסטים - במיוחד שכוללים ערכי Loss, גרדיאנטים קטנים או עדכוני משקלות - מבוצעים בדיקוק מלא .FP32

הפרדה זו מאפשרת לנצל את היתרון של מהירות זיכרון ב-FP16, מבלי לפגוע ביציבות המודל. המעבר מנוהל אוטומטית על ידי פרימיוםורקים AMP - Automatic Mixed Precision, במערכות כלים כגון TensorFlow ו-PyTorch .Precision

AMP ב-PyTorch: דוגמה פשוטה

הופך PyTorch את השימוש ב-Mixed Precision: לפניו עצם רכיבים מרכזיים autocast, torch.cuda.amp SMBCELoss, שמבצע המרות אוטומטיות ל-FP16 היכן שנitin ו-zeta, GradScaler, שמנגן מפני Underflow על ידי סקילינג של Loss לפניו Backpropagation ו- Unscaling. דוגמה מלאה לאימון ResNet18 על CIFAR-10:

```
import torch
from torch import nn, optim
from torchvision.models import resnet18
from torchvision.datasets import CIFAR10
from torchvision.transforms import ToTensor
```

```

from torch.utils.data import DataLoader
from torch.cuda.amp import autocast, GradScaler
from tqdm import tqdm

device = 'cuda' if torch.cuda.is_available() else 'cpu'

train_data = CIFAR10(root='.', train=True, download=True, transform=ToTensor
())
train_loader = DataLoader(train_data, batch_size=64, shuffle=True)

model = resnet18(num_classes=10).to(device)
optimizer = optim.Adam(model.parameters())
loss_fn = nn.CrossEntropyLoss()
scaler = GradScaler()

for epoch in range(3):
    total_loss = 0
    for x, y in tqdm(train_loader):
        x, y = x.to(device), y.to(device)
        optimizer.zero_grad()
        with autocast():
            pred = model(x)
            loss = loss_fn(pred, y)
            scaler.scale(loss).backward()
            scaler.step(optimizer)
            scaler.update()
            total_loss += loss.item()
    print(f"Epoch {epoch+1}, Loss: {total_loss:.3f}")

```

Listing 18.1 AMP: ResNet18 על CIFAR-10 עם אימון

לולאת האימון בגרסה זו דורשת רק שלושה שינויים: עטיפת ה-Forward scaler.scale(loss).backward().backward() loss.backward() Pass, ב-AMP, ב-החלפת () loss.backwards() בקריאה ל-(), scaler.step(optimizer) התאימות קטנות ושימוש ב- scaler.step() במקום עד האופטימיזציה הסטנדרטי.

אלו יכולות להביא להאוצאות דרמטיות - במיוחד באימון מודלים גדולים או בשימוש בבאץ'ים גדולים.

איך לדעת אם זה מועיל

כדי למדוד את התועלת האמיתית של Mixed Precision, התחילה במעקב אחרי זמן ריצה. ניתן להשתמש ב-`time` בתוכה תא אימון או להשוות אפוקים. לגבי זיכרון, הריצו `nvidia-smi` בתא נפרד וצפו בירידת ביצריה. אחד היתרונות המידיים הוא האפשרות להגדיל את גודל הבאץ' מבלתי להיות מוגבלים זיכרון - לעיתים להכפיל או לשולש את ה-`batch_size`.

מתי זה עובד ומתי להיזהר

Mixed Precision Training עובד היטב ברוב התרחישים - אך לא בכלם. מודלים עמוקים או רחבים מאוד, כמו Transformers, מרווחים במיוחד. גם שימושות עם קלטים גדולים, כמו יצירת תמונות או עיבוד רצפים. במקרים אלו, AMP צריכה להיחשב כברירת מחדל.

עם זאת, יש חריגים: ב-GPU ללא Tensor Cores - כמו קרטיסי Kepler או Maxwell ישנים - יתכן שלא תראו שיפור ואף אלה. במשימות שדורשות גילי אונמליות עדינות, כמו אבחונים רפואיים או סיווג אירועים נדירים - ירידת הדיקן הנומי **עלולה לפגוע** ברגישות המודל. לבסוף, ישקרים חריגים שבהם נוצרים NaNs באימון עקב Underflow או Exploding Gradients. במקרה זה ודאו שאתם משתמשים ב-`GradScaler` או הבעה נמשכת, הורידו את קצב הלמידה או חזרו זמניFP32.

ברירת מחדל חדשה ב-2025

נכון ל-2025, אкосיסטם הלמידה העמוקה מניח שימוש ב-Mixed Precision סביבות כמו TensorFlow PyTorch Lightning כבר מגדירותAMP כברירת מחדל. מרבית המחברות ב-Kaggle משתמשות בזה באופן מובנה. גם

בـ-Hugging Face הרבה מסקריפטי האימון מפעלים זאת כברירת מחדל. אלא אם אתם במחקר שדורש במפורש דיוק נומרי גבוה - אין עוד סיבה לא להשתמש בـAMP.

אם אתם מאמנים מודל גדול, משתמשים בـGPU מודרני, או עובדים עם באכ'ים מעל 128 - עלייכם להשתמש בـMixed Precision. היתרונות: התוכניות מהירה יותר ופחות צוואר בזיכרון זיברונו.

סיכום

Mixed Precision Training היא אחת מהטכניקות הנדרות בלמידה עמוקה שמשפרות גם ביצועים וגם יעילות - מבליל לדרוש שינוי ארכיטקטוני. היא מאפשרת למודלים לróż מהר יותר ולצורך פחות זיכרון ועדיין הגיעו לאותן תוצאות (או אפילו טובות יותר). בזכות כלים כמו AMP, שילוב בתהיליך האימון הפך לטריונייאלי.

הweeneyון הבסיסי פשוט: חשבו מהר היכן שניתן וחשבו בזיהירות היכן שחוובה. כמו בכל הנדסה טובה - מדובר על איזון. עידן של מודלים גדולים ותקציבים מוגבלים, האיזון הזה יכול לעשות את כל ההבדל.

שיעור 19

Fine-Tuning ו- Transfer Learning

מבוא

אימון מודל למידה عمוקה מאפס דורש לרוב דאטاهסטים עצומים, משאים המשובים משמעותיים וכיול נרחב. בפועל, לעיתים קרובות יעיל יותר לא להתחיל מאפס, אלא ממודל מאומן מראש - מודל שכבר למד "יכונים שימושיים מדעתה גдол יותר. גישה זו נקראת **Transfer Learning**.

Transfer Learning מנצל את יכולות חילוץ הפיצ'רים של מודל שאומן על דומיין כללי (למשל ImageNet או Wikipedia) ומתאים אותו לשימושה "יעודית" (למשל סיוג תמונות רפואיות או ניתוח סנטימנט). התהילה של התאמת המודל המאומן מראש לשימוש החדש נקרא **Fine-Tuning**.

בשיעור זה נציג את המוטיבציה, האסטרטגיות והיישום של Transfer Learning ו-**Fine-Tuning**, עם דוגמאות מתחום הראייה הממוחשבת (CV) ו-NLP.

מהו ?Transfer Learning

בלמידה מונחית קלאסית, אנו מאמנים מודל משקولات מאופסות אקראית בעזרת דוגמאות מתויות. ב-**Transfer Learning**, אנו מתחילהים ממודל שמשקولاتיו כבר אומנו על דאטהסט רחב ודומה. בראשות עמודות, השכבות המוקדמות נוטות ללמידה פיצ'רים כלליים

(कցוות, טקסטורות, דפוסים מקומיים), בעוד השכבות המאוחרות לומדות "יצוגים ספציפיים לדומיין". Transfer Learning מניה שייצוגים ברמה נמוכה אלה ניתנים לשימוש חוזר בין משימות שונות. דוגמה: רשות קונבולוציה (CNN) שאומנה לסיווג אובייקטים בתמונות טבע (למשל כלבים, משאיות) יכולה להיות מותאמת לסיווג תמונות מוצרים או צילומי רנטגן באמצעות שימוש חוזר בשכבות המוקדמות ועדכון רק של שכבות הסיווג הסופיות.

מודלים מאומנים נפוצים

- ראייה ממוחשבת (CV):** Transformers Vision EfficientNet, MobileNet, ResNet, .ImageNet - לרוב אומנו על (ViT)

- עיבוד שפה טבעיות (NLP):** DistilBERT GPT, RoBERTa, BERT, . - לרוב אומנו על קורפוסי טקסט גדולים.

- אודיו / דיבור:** Whisper wav2vec, - אומנו על גלי קול או ספקטrogramות.

מודלים מאומנים זמינים ברפואיטורי ציבוריים כמו Hugging Face Model או Torchvision Hub וניתן לטעון אותם בשורת קוד אחת.

שלוש אסטרטגיות Fine-Tuning

בדרכם כלל מיישמים Transfer Learning באחת שלוש דרכים:

- Feature Extraction (Fixed Base):** מ קופאים את כל השכבות המאומנות מראש ומאמנים רק שכבת פלט חדשה למשימה.

- Partial Fine-Tuning:** מ קופאים שכבות מוקדמות ומבצעים אימון לשכבות המאוחרות ולרأس הפלט.

- Full Fine-Tuning:** מאפשרים לכל המשקولات להתעדכן, לרוב עם קצב לימידה נמוך משמעותית.

כל גישה מאזנת בין גמישות לציבורות. Full Fine-Tuning עשוי להניב ביצועים גבוהים יותר במשימות ייחודיות לדומיין, אך דורש רגולרייזציה קפדנית.

דוגמה: MNIST Fashion ל- ResNet על Fine-Tuning

נניח שאנו רוצים לסוג תמונה בגווני אפור של בגדים מותך Fashion MNIST באמצעות Transfer Learning.

שלב 1: טעינת מודל מאומן מראש

```
from torchvision.models import resnet18
model = resnet18(pretrained=True)
```

שלב 2: הקפתה כל השכבות

```
for param in model.parameters():
    param.requires_grad = False
```

שלב 3: החלפת שכבת הפלט

```
import torch.nn as nn
model.fc = nn.Linear(512, 10)
```

שלב 4: אימון ראש הפלט

נשתמש בקצב למידה קטן (למשל 6×10^{-6}) ונאמן את שכבת הסיווג החדשה. ניתן בהמשך לשחרר בהדרגה שכבות נוספות.

שיקולים מתמטיים

במהלך Fine-Tuning, חשוב לשתמש בקצב למידה η מתאים. נסמן θ כוקטור הפרמטרים של המודל המאומן מראש ו- L כפונקציית המבחן למשימה החדשיה. עד העדכון הוא:

$$(19.1) \quad \theta_{k+1} = \theta_k - \eta \cdot \nabla_{\theta} L(\theta_k)$$

בדרכם כלל מיישמים קצב למידה קטן יותר לשכבות המאומנות מראש וקצב גבוה יותר לשכבות החדשיה. כך נמנע Catastrophic Forgetting של הייצוגים שנלמדו.

Fine-Tuning ב-NLP

מודלים לשפה טבעיות כמו BERT מותאמים בגישה דומה. למשל:

```
from transformers import BertForSequenceClassification

model = BertForSequenceClassification.from_pretrained(
    "bert-base-uncased", num_labels = 2 )
```

פקודה זו טוענת את BERT עם ראש סיווג מותאם לניתוח סנטימנט בינהרי. תהליך זה Fine-Tuning מתבצע באמצעות Backpropagation רגיל על טקסט מותpig.

מתי לא לשתמש ב-

למרות עצמתו, לש-Transfer Learning מגבלות:

- כאשר דומיין המקור והדומיין היעד שונים מאוד (למשל כלבים ← סריוקוט רפואיות).
- כאשר המודל המאומן מראש מכניס הטוויות שאינן רלוונטיות למשימת החדשיה.

- כאשר יש מספיק DATA ייודי לדומין ומشبבים חיישביים כדי לאמן מאפס.

תמיד יש להשוות מול בסיס של מודל מאומן מאפס.

טיפים

- התחילה תמיד עם `requires_grad = False` וחררו בהדרגה שכבות נוספות.
- השתמשו בקצב למדיה נמוך לשכבות המאומנות מראש.
- עקרו אחרי ביצועי הולידציה כדי לוחות Overfitting.
- במודלים גדולים השתמשו בקצבים למדיה שונים לשכבות (LR Layer-wise Scheduling).

סיכום

Transfer Learning מאפשר לנו לבנות מודלים מדויקים במהירות, עם פחות DATA ופחות משבבים. הוא פועל על ידי שימוש חוזר בייצוגים כלליים ממודל בסיס שאומן על משימה קרובה וההתאמתם באמצעות אמצעות Fine-Tuning. גישה זו היא כירום סטנדרט כמעט בכל תחומי של למדיה عمוקה. בין אם מדובר בסיווג תמונות או בניתוח טקסט, Transfer Learning משנה את תהליך פיתוח המודל מאימון להתאמת - מאיץ את הפיתוח וմשפר את ההכללה.

20 שיעור

שמירה, טעינה ו-*Versioning* של מודלים

מבוא

לאחר השקעת שעות - או ימים - באימון מודל למידה عمוקה, הדבר האחרון שתרצטו הוא לאבד את התוצאה. דיק גבוח אינו אומר דבר אם לא ניתן לעשות שימוש חוזר במודל, לשחזרו, או להטמיו.

שיעור זה עוסק בפרקטיקה מקצועית של שמירה, טעינה ו-*Versioning* של מודלים באופן שתומך בשיתוף פעולה, שחזריות ותחזוקה ארוכת-ពנים. נבחן את ההבדלים בין פורמטים כמו `SavedModel`, `TorchScript` ו-`ONNX` ונראה דוגמה מעשית עם CNN שאמן על `Fashion MNIST`.

חשיבות שמירת מודלים

שמירת מודל אינה רק צורך טכני - זהו עניין הנדי. אם איןכם שומרים את הפרמטרים שאומנו, מבנה המודל וסביבת העבודה שבה נוצר - למעשה אייבדתם את הניסוי. גם תוכאה של 94% דיק חסרת משמעות אם לא ניתן לשחזרה.

בעולם האמיתי, אי שמירת מודלים ותיעוד קונפיגורציות עלול להוביל לאיבוד נתונים, בזבוז זמן ופגיעה בעבודת המשך. שמירה נכון היא התחייבות עצמכם בעתיד - או לפחות הבא שייעבוד על הקוד שלכם.

שלושת פורמיות השמירה המרכזיות

באקויסיטטם של מידעה عمוקה ישנים שלושה פורמיות עיקריות לייצוא ושמירת מודלים מאומנים.

(PyTorch) TorchScript

הוא פתרון PyTorch לייצוא מודלים לפורמט עצמאי. הוא מקמל מודל לגרף חישוב סטטי, המאפשר רכזה מחוץ להרצת פיתון - למשל באפליקציות C++ או בموבייל.

```
traced = torch.jit.trace(model, example_input)
traced.save("model.pt")
```

Listing :20.1 שמירה TorchScript

לטעינה ושימוש ב:

```
model = torch.jit.load("model.pt")
model.eval()
```

Listing :20.2 טעינה TorchScript

הפורמט קומפקטי, מהיר ונitin להעברה - אידיאלי לאינטגרציה בפרויקט עם תלות מינימלית.

(TensorFlow) SavedModel

ב-TensorFlow, פורמט ברירת המחדל הוא SavedModel. הוא שומר לא רק את המשקלות וגרף החישוב, אלא גם חתימות פונקציה ומטא-דאטा. השמירה פשוטה:

```
model.save("my_model")
```

Listing :20.3 שמירה TensorFlow

וטעינה פשוטה באותה מידעה:

```
loaded = tf.keras.models.load_model("my_model")
```

Listing :20.4 טעינת SavedModel

פורמט זה מתאים במיוחד למימוש דרך TensorFlow ,TensorFlow Serving או Google Cloud AI Platform ,Lite .

(Open Neural Network Exchange) ONNX

ONNX הוא פורטט פתוח חזחה-פרוייקטים. מודל שאותן ב PyTorch יכול להיות מוצא ל-ONNX ולאחר מכן להיות ממומש בפרויקטים שונים או שפות אחרות.

```
torch.onnx.export(model, dummy_input, "model.onnx",
                  input_names=["input"], output_names=["output"])
```

Listing :20.5 PyTorch ל-ONNX ייצוא

ONNX אידיאלי לפירשת מודלים בסביבות שאין קשורות לספרייה מסוימת, כגון REST APIs Java, C++ או Triton In-ference Server .ONNX Runtime

דוגמה : CNN על Fashion MNIST

נבחן דוגמה קונקרטית: נגיד ונאמן רשות קונבנצייה קטנה לסיווג תמונות ונשמר אותה בשלושה פורמטים: TorchScript, משקولات גולמיות, ONNX ו-DNNX .

```
import torch
import torch.nn as nn

class Net(nn.Module):
    def __init__(self):
        super().__init__()
        self.net = nn.Sequential(
            nn.Conv2d(1, 16, 3, padding= 1),
```

```

        nn.ReLU(),
        nn.Flatten(),
        nn.Linear(16 * 28 * 28, 10)
    )
def forward(self, x) :
    return self.net(x)

model = Net()

```

Listing :20.6 CNN ב- PyTorch

ראשית, ניתן לשמר רק את המשקלות:

```

torch.save(model.state_dict(), "cnn_weights.pth")

model.load_state_dict(torch.load("cnn_weights.pth"))
model.eval()

```

Listing :20.7 שמירה+טעינה state dict

לאחר מכן, לימוש מחוץ ל- PyTorch :

```

example_input = torch.randn(1, 1, 28, )28
traced = torch.jit.trace(model, example_input)
traced.save("cnn_scripted.pt")

```

Listing :20.8 שמירה TorchScript

ולבסוף, לייצוא ל-ONNX :

```

dummy_input = torch.randn(1, 1, 28, )28
torch.onnx.export(model, dummy_input, "cnn.onnx",
                  input_names[="input"], output_names[="output"])

```

Listing :20.9 ייצוא ONNX ל-

בחירה הפורמט המתאים

הפורמט המתאים תלוי בהקשר השימוש, באקוסיסטם ובמטרות. לכל אפשרות יתרונות ייחודיים, כפי שמסוכם בטבלה:

פורמט	שימוש עיקרי	יתרונות
TorchScript	איינפנסיס מהיר, איןטרציה ל-++/C+/mobilly	קל משקל, נייד, רץ ללא פיצ'ון
SavedModel	Serving TensorFlow או TFLite	כולל גרפ, משקלות וחתימות; ידידותי לענן
ONNX	מעבר פרימוררים בין רבים (Triton, ONNX Runtime, Optimum)	פורמט פתוח, נתמך ע"י רנטיטים

טבלה 20.1: השוואת פורמטי ייצוא מודלים

Versioning

שמירת מודל היא רק חלק מהסיפור. פרויקטים רציניים דורשים גם Versioning. כל מודל שנשמר צריך להיות עם שם גרסה ברור, למשל model_v2.1.onnx. בנוסף, מומלץ לכלולקובץ לוג או מטא-דאטה נלווה, המכיל: תאריך אימון, תיאור קצר של השינויים, היפר-פרמטרים מרכזיים ומדדי ביצוע סופיים על DATAHSET הטעט. מידע זה מבטיח שאתה - או כל אדם אחר - תוכלו לשחזר את התוצאות גם חודשים מאוחר יותר. שקלו לשלב כלים כמו Git, MLflow, או DVC לניהול ניסויים וארטיפקטים של מודלים לאורך זמן ובין מכונות.

סיכום

שמירת מודל אינה סוף האימון - אלא תחילת חייו בעולם האמיתי. בין אם מדובר במימוש לפודקשן, שיתוף עם הצוות, או חזרה לניסוי ישן - סרייאלייזציה נכונה ו- Versioning חיוניים.

פורמטים כמו TorchScript, SavedModel ו-ONNX מאפשרים להעביר מודלים בין שפות, רנטיים ותשתיות. אך אף פורמט אינו יכול לשמר את הכוונה או ההקשר - זו כבר אחריותם.

21 שיעור

IMPLEMENTATION

INTRODUCTION

לאחר אימון ושמירת מודל למידה عمוקה, השלב החינוי הבא הוא להפוך אותו לשימושי - לא רק מותוך המחברת שלכם, אלא גם על ידי מערכות אחרות. זהו הקשר בין מודולינג לפרודקشن. כך מסווג תומונות הופך לשירות בזמן אמיתי, כך מנוע המלצות מניע אפליקטיבית ווב וכך למידה عمוקה הופכת לחלק ממוצר עובד.

בשיעור זה נראה כיצד לפרסום מודל סיווג תומונות מאומן כ-`API RESTful` באמצעות `TorchScript` ו-`FastAPI`. המטרה: לקבל תמונה דרך `HTTP`, להציג תחזית בפורמט `JSON` - נקי, מהיר ומוכן לפרודקشن.

MODEL FOR SERVICES

בסביבת פרודקشن לא מרכיבים מחדש סקריפטי אימון או מחברות כדי לקבל תחזיות. במקום זאת, עוטפים את המודל בשירותות מותמץ - זהה שמאזין לבקשות כניסה, מעבד אותן ומחזיר תחזיות. ארכיטקטורה זו מאפשרת לשלב את המודל בארכי אינטרנט, אפליקציות מובייל, דאשборדים, מערכות `IoT`, או כל ממשק תוכנה אחר.

API של מודל מקבל קלט (למשל תמונה), מבצע פרא-פרוססינג, מפעיל אינפראנס באמצעות מודל טעון מראש ומחזיר את התוצאה במבנה מסודר

כמו JSON.

שירות אינפראנס עם בהירות FastAPI

FastAPI הוא פרויקט פיתוח לבניית APIs עתירי ביצועים. הוא משתמש בטבעיות עם NumPy, PIL, PyTorch ותומך בעיבוד אסינכרוני ולידציה אוטומטית של בקשות ותיעוד אינטראקטיבי בזמן אמת דרך Swagger. העיצוב שלו הופך אותו מותאים במיוחד לשירותי אינפראנס: קל משקל, מהיר לעלות ופשוט לבדיקה.

סיווג תמונות עם Case Use TorchScript

נניח שאימנו רשת קומבולוציה על Fashion MNIST, "צאננו אותה כ-cnn_scripted.pt" וANO רוצים לחושף אותה כ-API REST שמקבלת תמונה ומחזירה את המחלקה החזויה. להלן המימוש המלא:

```
from fastapi import FastAPI, File, UploadFile
from fastapi.responses import JSONResponse
from PIL import Image
import torch
import torchvision.transforms as transforms
from io import BytesIO

# Device configuration
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

# Class labels for Fashion MNIST
LABELS = {
    :0 "T-shirt/top", :1 "Trouser", :2 "Pullover",
    :3 "Dress",      :4 "Coat",     :5 "Sandal",
    :6 "Shirt",       :7 "Sneaker",   :8 "Bag",     :9 "Ankle boot"
}
```

```

# Load model once on startup
model = torch.jit.load("cnn_scripted.pt", map_location=device)
model.eval()

# Define image preprocessing
transform = transforms.Compose([
    transforms.Grayscale(),
    transforms.Resize((28, 28)),
    transforms.ToTensor()
])

app = FastAPI()

@app.post("/predict/")
async def predict(file: UploadFile = File(...)):

    try:
        image = Image.open(BytesIO(await file.read())).convert("RGB")
        x = transform(image).unsqueeze(0).to(device)

        with torch.no_grad():
            output = model(x)
            probs = torch.softmax(output, dim=1)=
            pred_idx = probs.argmax(dim.)1=item()
            confidence = probs.max(dim.)1=item()
            label = LABELS[pred_idx]

        return JSONResponse(content={
            "prediction_index": pred_idx,
            "label": label,
            "confidence": round(confidence, 4)
        })

    except Exception as e:
        return JSONResponse(status_code=400, content={"error": str(e)})

```

Listing :21.1 FastAPI+TorchScript

השירות טוען את המודל פעם אחת ליזכרון, מבצע פרה-פרוססינג לתמונה שהועלתה כך שתתאים לקלט המודל ומחזיר את האינדקס של המחלקה עם תווית וביתחון התחזית.

הרצה השירות מקומי

כדי להריץ את השירות מקומי, שמרו את הקוד בקובץ api.py ואז הריצו:

```
unicorn api:app --host 0.0.0.0 --port 8000 --reload
```

Listing :21.2 FastAPI עם Unicorn

לאחר ההרצה, בקרו בכתובת:

<http://127.0.0.1:8000/docs>

FastAPI תיצור באופן אוטומטי ממשק web שבו תוכלו להעלות תמונות, להריץ תחזיות ולבדוק תוצאות. זו אחת הסיבות לפופולריות שלו בקרב מתכננים בתחום - בדיקת מודל הופכת פשוטה כמו שימוש באתר.

מעבר לענן

לאחר שהשירות פועל מקומי, ניתן למעשה אותו בכל סביבה מודרנית. אפשרויות נפוצות כוללות:

- אחסון בפלטפורמות כמו Render או Hugging Face
- הרצה ב-Colab עם ngrok לחשיפת פורט מקומי
- מימוש במכונה וירטואלית ב-GCP, AWS, או Azure, לרוב בתוך קונטינר Docker

בכל המקרים, המבנה נשאר זהה, אך נקודת הכניסה היא כעת Endpoint ציבורי שMOVED להציג תחזיות.

טיפים ליציבות וביצועים

כדי להבטיח שה-API שלכם יתנהג בעקביות ובייעילות בפודקשן:

- הקפידו להגדיר `model.eval()` לפני אינפראנס כדי להבטיח התנהגות דטרמיניסטית.
- עטפו את בлок התחזית ב-`torch.no_grad()` כדי לחסוך בזכרון ולנטרל מעקב גרדיאנטים.
- הימנעו מטעינת המודל בכל בקשה - טעינו אותו פעם אחת כשהשרת עליה.
- השתמשו ב-`map_location=torch.device("cpu")` בטעינה מודלים בסביבות GPU אחרת, יתכונו שגיאות אם המודל נשמר על GPU.
- לשיפור מהירותת תגובהה, שמרו את טרנספורמציות התמונות קלות ובצעו עיבוד באותו התקן שבו רץ המודל.
- אם אתם מתכונים לטפל במספר בקשות לשניה, בטלו את דגל `-reload` של FastAPI והשתמשו בשרת פודקשן כמו Gunicorn או Uvicorn עם מספר עובדים. במערכות מתקדמות ניתן למש用工 Batch Inference או תורים אסינכרוניים באמצעות Celery או Ray Serve.

סיכום

שימוש מודל למידה עמוקה הוא הרגע שבו המודל הופך למערכת. API פשוט, המבוסס על API FastAPI ו-TorchScript, יכול להפוך סקריפט מקומי לשירות רב-שימושי, ניתן לבדיקה ונגיש.

לא נדרש לכך הרבה קוד. אך כן נדרש תכנון ברור של משקדים, התקנים ויציבות. מרגע שהמודל מומש, הוא כבר לא רק קובץ - אלא חלק ממוצר.

שיעור 22

התמחויות

מבוא

אם הגיעם لنקודה זו, כבר יש לכם בסיס טוב בלמידה عمוקה. אתם מבינים כיצד רשותות ניירונים נבנות, מתאמנות, מאופטומות ומיושמות. אתם יודעים מה זה אומר לספק תחזיות בזמן אמת וכי怎ן נוהל מודלים באופן מבודק ובהירות.

אך במידה עמוקה איןנה תחום יחיד - היא אкосיסטם של תחומיים متמחים. כל תחום מגע עם האתגרים שלו, ארכיטקטורות מודלים ייחודיות ומוסכימות משלו. בשיעור זה נמפה את ארבעת הcientists הנפוצים ביותר להתחמות מקצועית: ראייה ממוחשבת, עיבוד שפה טבעית, למידה טבלאית וניתוח סדרות זמן.

מדוע התמחות חשובה

אין ארכיטקטורת מודל אחת שמתאימה לכל הבעיות. מודל שמצטיין בזיהוי תסבוכות עלול להיכשל לחלוטין כאשר מיושם על טקסט. שיטה שפותחה לדאטה טבלאי של CRM העשויה להתקשות ללמידה תלות טמפורלית בתחזיות פיננסיות. התמחות בלמידה עמוקה איננה בחרית אלגוריתם מועדף - אלא התאמת הכלים לבנייה הדאטה ולطبע המשימה. הבנת המיפוי הזה היא המפתח לעובדה עיליה בלמידה عمוקה יישומית.

ראייה ממוחשבת (Computer Vision)

дана и израильская индустрия изгibt образованием «иудеев». Благодаря локаторам пространственных структур, такие как CNNs, модели могут строить структуру изображения, определять местоположение объектов и их взаимное расположение.

Решетки RNN включают в себя не только пространственные структуры, но и временные, что позволяет им обрабатывать последовательности изображений. Для этого используются различные архитектуры, такие как LSTM и GRU, которые позволяют сохранять информацию о предыдущих кадрах и использовать ее для предсказания будущих. Благодаря этому, модели могут выполнять задачи, такие как классификация изображений, распознавание лиц, сегментация изображений и даже генерация изображений на основе текстовых описаний.

Большинство моделей RNN, таких как LSTM и GRU, являются рекуррентными, то есть они обрабатывают последовательность изображений, один за другим. Это ограничивает их способность обрабатывать изображения в реальном времени, так как они должны ждать, пока предыдущий кадр будет обработан, чтобы продолжить обработку следующего. Для решения этой проблемы были разработаны альтернативные архитектуры, такие как Transformer, которые позволяют обрабатывать изображения одновременно во всех точках, что делает их гораздо более эффективными для задач обработки изображений.

מודלי שפה (Language Models)

дана текстуальный он правдив, тесно связан и имеет семантическую природу. Благодаря этому, он может быть легко расширять и комбинировать с другими текстами.

архитектура центральная, которая преобразует текст в векторные представления, которые затем используются для выполнения различных задач, таких как классификация, распознавание языка и генерация текста. Модели LSTMs и GRUs, которые мы рассмотрели ранее, являются примерами альтернативных архитектур, которые также могут использоваться для выполнения этих задач.

модель способна обучаться на больших объемах текста и использовать полученные знания для выполнения различных задач. Одним из основных преимуществ является то, что модель может обрабатывать текст в реальном времени, что делает ее особенно полезной для задач, требующих быстрой обработки информации.

использование альтернативных архитектур, таких как Transformer, позволяет обрабатывать текст одновременно во всех точках, что делает его гораздо более эффективным для задач обработки текста.

Минусом является то, что эти модели требуют очень много вычислительных ресурсов и памяти для обучения и хранения больших объемов данных.

מייל ששלח או מייצרת סיקום קצר של מסמך משפטי. המודל נדרש להבין, לתמצעת ולפעול - הכל מקלט טקסטואלי בלתי מובנה.

דатаה טבלאי (Tabular Data)

דータה טבלאי מצוי בכל תחום עסק: יומני מכירות, יצוא CRM, רישומי פיננסים ומטא-דータה על לקוחות. למרות שהפורמט נראה פשוט - שורות ועמודות - האתגרים המודולאים אמתיים. פיצרים קטגוריים, ערכאים חסריים, מחלקות לא מאוזנות וمزוהים בעלי קרדינליות גבוהה מסובכים את הלמידה. מודלים מסורתיים מבוסטי Ensembles כמו XGBoost ו-LightGBM שולטים בתחרויות טבלאיות רבות. הם מהירים, רובוטיים וניתנים לפרשנות. עם זאת, גישות נוירוניות מפותחות. MLPs עובדים היטב כאשר יש כמות דータה גדולה וארכיטקטורות כמו TabNet או FT-Transformer משלבות מנגןוני Attention כדי למודל יחסים בין עמודות בקרה עיליה יותר. בתחומים כמו דירוג אשראי או חיזוי נתיחה, פרשנות היא חובה. מגבלות רגולטוריות ותפעוליות דורשות לעתים להצדיק מדוע התקבלה תחזית - למשל באמצעות ערכי SHAP או שיטות מבוססות Permutation.

סדרות זמן ותחזיות במידע רציף (Time Series)

דータה טמפורלי מוסיף ציר חדש: הזמן. כל תחזית תלואה לא רק בפיצ'רים, אלא גם במה שקדם לה. חיזוי סדרות זמן כולל לרוב מגמות, עונתיות, עיכובים או ערי מדידה - ודורש מודלים שלומדים גם רצף וגם הקשר. לשם כך משתמשים במודלים כמו LSTM ו-GRU, השומרים זיכרון פנימי לאורץ זמן. מודלים מבוססי CNN כמו TCN (Temporal Convolutional Networks) מסננים על חלונות זמן ומספקים אלטרנטיבה לרקורסיה. לאחרונה וריאנטים של Transformer כמו TFT או Autoformer מציעים חלופות סקלאbialיות למדול רציפים ארוכים. עברו אותן בתדריות גבוהה, כמו דטה מחישנים או מערכות בקרה, מודלי State Space - כגון Mamba - צוברים פופולריות. דמיינו חיזוי ביקוש למוצר על פני מאות קודי מוצר SKU, כל אחד עם

דפוס משלו, היסטוריית מביצעים וסוגות גיאוגרפיה. המטרה איננה רק לאזכור את העבר - אלא ללמידה כיצד הוא מתפתח ומתי הוא משתמש.

כיצד לבחור תחום להמשך

לכל תחום יש שיקולים שונים - לא רק במבנה הדאטה, אלא גם במטרות המודול, מגבלות האינפראנס וaintegration מערכתי. אין פתרון אחד שמתאים לכלו.

כלל אצבע פשוט:

- אם הקלט שלכם הוא תמונה - התחילה ממודל קונבולוציה.
 - אם אתם עובדים עם טקסט חופשי - בדקו *Transformers*.
 - אם יש לכם שורות ועמודות - שקלו מודלי *Ensemble* או *MLPs*.
 - אם הדאטה שלכם מותפת לאורך זמן - חקרו ארכיטקטורות רציפות מבוססות זמן או מבוססות *Transformer*.
- כרגע, בניית הדאטה צריך להנחות את הבחירה - לא הפופולריות של האלגוריתם.

סיכום

למידה עמוקה היא משפחה של גישות המותאמות לסוגי דאטה שונים ונתגרים מודולים מגוונים. ראייה, שפה, טבלאות ורכפים - כל אחד דורש דרך חשיבה אחרת. בחירת תחום התמחות תלוי בדומין, במטרות ובטבע הבעיה. יש מהנדסים עמוקים במודאליות אחת. אחרים לומדים לשלב ביניהן.

בחרו את הכלים לפי צורת הדאטה שלכם - לא לפי האופנה בתחום.

23 שיעור

מפתח דרכים למתודת DL בתעשייה

מבוא

ברכוות - הגיעם לשיעור האחרון בקורס יסודות הלמידה העמוקה המודרנית. עברנו דרך ארוכה: מהעברית הפנימית של רשות ניירונים ועד מימוש שירות תחזיות בזמן אמיתי בענן.

עת הגיאו הזמן לעזר ולשאול שאלה עמוקה יותר: כיצד נראה הדרך קדימה עبور מהנדס למדידה עמוקה בעולם האמיתי? שיעור זה אינו עוסק בקוד חדש או במודלים חדשים. הוא עוסק בגישה מנטלית. במבנה מערכות שמחזיקות לאורך זמן, בקבלת החלטות שניתן לעמוד מאחוריהן ובഫיכת עבודה מניסוי להשפעה ממשית.

מודול למערכת פרודקشن

aimon מודול מודיק הוא רק התחלה. בסביבות תעשייתית, דיק בלבד אינו מספיק. מודול טוב הוא מודול שנייתן למש, לנטר, לעדכן ולסמן עליו לאורך זמן.

משמעות הדבר היא שינוי חשיבה מ- "המודול" ל- "המערכת". קוובץ עם משקלות איננו פתרון. הטמעה מלאה - מאימון וolidציה ועד שירות וניתוח - הוא זה שהופך למדידה עמוקה לאופרטיבית באמת.

שלושת רכיבי הליבה של מערכת DL

לכל מערכת DL בעולם האמתי יש שלושהרכיבים יסודיים:

- הראשון הוא רכיב האימון. כולל דאטה, קונפיגורציית היפר-פרמטרים, ארכיטקטורת מודל ויכולת שחזור ביצועים. תהליכי אימון מתוכנן היטב מאפשר לחזור על ניסויים חוזרים מאוחר יותר - עם אותן תוצאות.
- השני הוא רכיב הולידציה והניטור. זה מה שמספר אם המודל עדיין עובד טוב על דאטה חדשה, אם יש ירידה בביטויים ואם נדרש פעולה. ללא ניטור מתמשך, מודל שהיה טוב עלול להפוך למסוכן בשקט.
- השלישי הוא הממשק - לרוב API או שירות סטרימינג - דרכו צורכים את התוצאות. זה מה שמערכות אחרות רואות. עליו להיות מהיר, יציב ושקוף.

שלושת החלקים הללו חייבים להיות מותאימים וմבוקרים. יחד הם מהווים את היסוד של כל מערכת ML בפודקשן.

מקרה לדוגמה: סיווג כוונת אימיילים

דמיינו שאתם נדרשים לבנות מודל שמצויה כוונת משתמש מתוך אימיילים נכנסים מלקוחות:

- מודול אימון שורום לוג לכל ניסוי, כולל:
 - גרסת הדאטאסט
 - היפר-פרמטרים
 - מדדים סופיים
- מודול ולידציה ש:
 - בודק את המודל مدى שבוע על טיקטי תמייה חדשים
 - מתריע על ירידה בדיק או בכיסוי
- API שמחזיר:

- תחזיות
- ציוני ביטחון
- תיאור טקסטואלי קצר

כל זה מנוטר באמצעות MLflow או מערכת מטא-דאטा מבוססת Git, כך שכל שינוי מתועד וניתן להחזירה. זהו הסטנדרט למידה عمוקה ברמה פרודקסן.

טיפים למתודס בעולם האמיתי

לעבוד באופן מקצועי בלמידה عمוקה זה לא רק לכתוב קוד טוב. נדרשת ממשמעת. תמיד תעמדו את מה שאתם עושים. לכל גרסה מודל חייב להיות תאריך ברור, קונפיגורציה ותוצאות הערכה. ודאו שחזוריות. מישחו - לעיתים אתם עצמכם - יctrיך להריץ את אותו אימון בעוד חצי שנה ולקבל אותה תוצאה. השתמשו בפורמיティ סריאלייזציה רובייטיים כמו TorchScript, ONNX או SavedModel. הימנוו מסקרים פטימי מותאמים אישית או שמירות אקרניות. לעולם אל תדלו על ולידציה. מודל שמתנהג בחוסר עקביות, או שמשנה התנחות לאורך זמן, אינו אמין. תכננו תהליכיים מודולריים. ה-UI שלכם לא צריך לדאוג אם המודל מתחתתיו השתנה. לוגיקת האימון צריכה להיות מבודדת מלוגיקת השירות. ותמיד הערכו עקביות לאורך זמן - לא רק דיווק חד-פעמי.

אתיקה ואחריות

למידה عمוקה יכולה להשפיע על גיבוב עובדים, אבחון רפואי, תמחור ועוד. עם כוח מגעה אחריות. לעולם אל תמשו מודל שלא נבדק על פני אוכלוסיות שונות של נתונים. תמיד הבינו מאיפה הגיע הדאטה - ומה יתכן שיחסר בו. היו מוכנים להסביר תחזיות, גם אם רק חלקית. והכי חשוב: שאלו את עצמכם האם הייתם סומכים על המערכת הזו אילו היא הייתה משפיעה

עליכם או על מישחו שייקר לכם. AI Responsible איןנו צ'קבוקס חוקי, אלא ערך הנדסי.

מקצוענות אמיתית

מה שմבדיל מהנדס מוחובבן אינו רק עומק המודלים שלו - אלא עומק החשיבה שלו. מהנדס DL מקצועי לא רק מכון היפר-פרמטרים. הוא בונה מערכות שהן עמידות, ניתנות לשחזר, ניתנות לניטור ובטוחות. זה דורש תשומת לב לפרטי זהה מה שהופך מודל מבטיח לפתרון אמיתי.

מניפסט למהנדס DL

העקרונות הבאים מסכימים את הערcis והגישה של מהנדס DL מקצועי בעולם האמתי. אלו אינם חוקים - אלא הרגלים שנבדקו וմבדילים מערכות בוגרות מפרטוטיפיים שביריריים:

- **כל מודל הוא מערכת.** המודל טוב רק כמו הדאטה, תהליך הולידציה וסבירת המימוש שבה הוא חי.
- **שחזריות מעל אינטואיציה.** כל ניסוי חייב להיות ניתן למעקב. מה שלא ניתן לשחזר - לא ניתן לสมוך עליו.
- **גרסאות.** דאטה, קוד, מודלים, קונפיגורציות, מדדים - עקבו אחרי כולם. קיצור הדרך של היום הוא הבאג של המחר.
- **ולידציה מעבר לדיווק.** הערכו גם Drift, יציבות, רוביוטיות וביתחון - רק מدد ייחיד על דאטاهסט סטטי.
- **תכננו לבשל.** ניטור בלבד אינו אופציה. הניחו שדברים לא יעבדו חלק. בנו התראות, ניסיונות חוזרים ואסטרטגיות גיבוי.
- **פרשנות.** אם מודל מקבל החלטה שאינכם יכולים להסביר - לא בטוח сможете לו לקבל החלטות בעולם האמתי.

- בחרו כלים לפי הדאטה - לא לפי טרנד. תנו לצורת הבעה להנחות את הארכיטקטורה - לא למה ש"חס" ב-[arXiv](#).
 - Bias הוא באג. חפשו אותו באופן אקטיבי. טיפול בו מוקדם. תעדו את דרכי ההתחמodoreות שלכם.
 - בנו אמונה. העבודה שלכם היא לא רק לאפטם מדדים - אלא ליצור מערכות שאחרים יכולים לסייע עליהם.
- אלו הם העקרונות שילו אתכם כשתעברו מבניית מודלים להנדסת מערכות אמינות. שמרו אותם.

לאן ממשיכים מכאן

הקורס הזה נתן לכם את היסודות. עכשו יש לכם את המונחים, הכלים והביטחון להתחילה פרויקטים אמיתיים. מכאן, תוכלו להעמק: לראייה ממוחשבת, מודלי שפה, במידה טבלאית או סדרות זמן. תוכלו להתמחות - או לפחות. תוכלו לעבוד בלבד - או בצוות. אך בכל דרך שתבחרו, אתם כבר מצוידים בצורת החשיבה הנכונה: שיטתיות, אתית ויישומית.

סיכום

קרירהה בלימדה عمוקה מותפתחת כל הזמן. אך עקרונות ההנדסה הנכונה נשאים: ללמידה, לבנות, למודד, לשתף. בחרו פרויקט קטן ואמיתי. התיחסו אליו ברצינות. תעדו את ההנחות. ממשו את הרעיון. אני מקווה שהקורס נתן לכם לא רק ידע - אלא גם ביטחון, כלים וモטיבציה להמשיך. העולם זוקק למערכות שתבנו - ואתם מוכנים.

בהצלחה!