"הסבר לתרגיל 3 בקורס "למידת מכונה"

נכתב על ידי תומר גיל עבור הקורס 89-511 באוניברסיטת בר-אילן שמעביר ד"ר יוסי קשת 9/5/2018

תודות מיוחדות: תמיר בוסקילה וטליה בן סימון

Contents

2	הקדמה
3	
3	טעינת המידע
3	מה יש במידע?
3	מה כל dataset אומר?
4	יצירת פרמטרים ובחירת היפר-פרמטרים
4	פרמטרים
4	היפר-פרמטרים
5	פסאודו-קוד
5	אימון הרשת
6	– Forward Pass – חישוב תוצאת הרשת
6	חישוב Loss
6	— Back Propagation – חישוב הגרדיאנטים
7	עדכון המשקלים
8	חישוב דיוק ו-loss על ה-dev_set
8	
8	Predicting on the Test Set
Q	רעיות פתרונות ונויפים

הקדמה

ראיתי שיש להרבה אנשים בעיה עם התרגיל, אז כתבתי הסבר "קצר" על מה עושים בתרגיל, או לפחות המסגרת של מה אמורים לעשות.

אני לא לוקח שום אחריות על איך התרגיל שלכם יצא אם תעבדו לפי המסמך הזה. השימוש על אחריותכם בלבד.

השתדלתי לחלק את התרגיל לחלקים לוגיים שונים, כדי שיהיה ברור מה צריך לעשות בכל חלק ואיך הם מתחברים. לא כל כך נכנסתי לתיאוריה ול"למה עושים מה שעושים". לא שזה לא חשוב, אבל זה לא קריטי לתרגיל (למבחן לעומת זאת...).

אני ממליץ מאוד לפני שמתחילים לעבור את ה-tutorial הזה שמלמד פייתון ויותר חשוב – שימוש ב- (SciPy אין צורך ב-(אין צורך ב-(SciPy) בכל מקרה מקווה שתיהנו מהתרגיל (או לפחות לא תסבלו) והמון בהצלחה!

תיאור התרגיל

אתחיל בלתאר מה בערך צריך לעשות בתרגיל.

טעינת המידע

החלק הראשוני של התוכנית היא קריאת הקבצים והפיכתם למטריצות של numpy. למזלנו vnp.loadtext(file_name). נשים לב שאחרי שהוצאנו את הקבצים יכול לעשות את זה בשבילנו עם np.loadtext(file_name). נשים לב שאחרי שהוצאנו את הקבצים מהיד rain_x מהפיאצה נקבל 3 קבצים: train_x (מכיל את מטריצת הדוגמאות עבור סט האימון ovalidation – כל שורה היא דוגמא מגודל 784), train_y (מכיל את התגים של הדוגמאות מהקובץ הנ"ל – כל שורה היא התיוג של השורה המתאימה מxim_x ו-stain_x (מכיל את הדוגמאות שעליהן הנ"ל – כל שורה היא התיוג של השורה המתאימה מxim_x (loadtext() (האמת נבדוק את המודל שלנו בסוף). נרצה לטעון את שלושתם למטריצות באמצעות (loadtext() (האמת שיש עוד קובץ בשם test.pred, אבל אנחנו נתעלם ממנו)

לאחר שנטען את הדאטה, נרצה ליצור מ-train_y ו-train_y שני סטי אימון: train ו-validation. בשביל זה אנחנו נעשה שני דברים: קודם נעשה shuffle להם¹ (או עם numpy של בשביל זה אנחנו נעשה שני דברים: קודם נעשה shuffle להם¹ (או עם train של סדר השורות פייתון – העיקר לוודא שכשאנחנו משנים את הסדר השורות ב-X לשנות בהתאמה את סדר השורות ב-Y), ואז נחלק אותו לשניים: train עם 80% מהמידע ו-dev

:פסאודו-קוד בסיסי

```
train x = np.loadtxt("train x")
train y = np.loadtxt("train y")
test x = np.loadtxt("test x")
shuffle(train_x, train_y)
dev_size = train_size * 0.2
dev_x, dev_y = train_x[-dev_size:, :], train_y[-dev_size:]
train_x, train_y = train_x[:-dev_size, :], train_y[:-dev_size]
```

מה יש במידע?

ה-dataset נקרא Fashion-MNIST שמכיל תמונות שחור-אפור-לבן בגודל 28*28 פיקסלים. כל פיקסל הוא float נקרא Pashion-MNIST שמכיל תמונות שחור-אפור-לבן בגודל 28*28 פיקסלים. כל פיקסל הוא float בטווח 0-255 שמתאר כמה כהה הפיקסל. (זה ה-x-ים) יש 10 מחלקות שהן סוגי לבוש וכד' שונים (יש רשימה בתרגיל) ולכל אחת יש אינדקס ייחודי בטווח בין 0-9 – זה בעצם מתאר לנו את התא בווקטור התוצאות של הרשת שאומר את ההסתברות שהתמונה היא מסוג המחלקה הזו. (זה ה-y-ים)

האם אכפת לכם? לא בטוח. בסופו של דבר כל x הוא וקטור בגודל 784 (28*28) (ראו הערה x האם אכפת לכם? לא בטוח. בסופו של דבר כל x הוא מספר בטווח 9-0. אתם תתאמנו עליהם בלי בהכרח לדעת איך נראית כל תמונה.

?מה כל dataset אומר

קיבלנו 3, שלכל אחד יש תפקיד שונה מבחינתנו בתהליך האימון:

- Train set הסט הגדול ביותר, זה הסט היחידי שהמודל שלנו יתאמן עליו (כלומר ילמד אותו). במהלך האימון נעבור על כל הtrain set כמה פעמים ובכל שלב נכניס אותו דרך הרשת, נבדוק כמה טעינו (loss), נחשב את הגרדיאנטים של כל פרמטר לפי ה-loss ונעדכן את המשקולות לפי הכלל של SGD (פירוט על כל זה בהמשך)
- loss מקבצי ה_train. עליו נרצה לבדוק train זה הסט המשלים של החולם. עליו נרצה לבדוק train זה הסט המשלים של החול לא פעם לא "ראה" (לא למד) ממוצע ודיוק עליו לכל אפוק. הוא בעצם דוגמאות שהמודל אף פעם לא "ראה" (לא למד) ואנחנו נבחן כמה טוב הוא בלסווג מידע כזה אתו. יעזור לנו לבחור היפר-פרמטרים.

מסודר train_x יהיה מאותו סוג (נגיד אם validation set- למה לעשות?shuffle הרנדומליות תדאג שלא כל ה-לפי סוג (נגיד אם clidation set) ובעצם תעזור לנו לכוונן את המודל למקרה כללי (דוגמאות שהמודל לא ראה מעולם יסווגו טוב)

Test set – נועד לתת לצוות הקורס לבחון אותנו (בתרגיל, בעולם האמתי נותן לעשות מעין מבחן סופי). אין לנו את התיוג שלו (מתעלמים מtest.pred) ואנחנו בעצם ניצור test.pred משלנו שמכיל את התיוגים שהמודל שלנו נתן לאחר שאימנו אותו על ה-train set.
 שימו לב! הציון יינתן על פי הדיוק של הקובץ הזה, אז אנחנו נרצה לתת את הזה שיצא מהמודל הטוב ביותר שלנו.

יצירת פרמטרים ובחירת היפר-פרמטרים

פרמטרים הם אותם מטריצות/ווקטורים משקולות שאנחנו נעדכן במהלך אימון המודל כדי לתת סיווג מדויק יותר – הם החלקים שהמודל ישנה לבד במהלך האימון.

היפר-פרמטרים הם אותם "פרמטרים" לתוכנית שאנחנו מחליטים עליהם, ויש להם השפעה עצומה על הצלחת המודל. המודל לא יכול ללמוד אותם כי הם חלק ממנו, כביכול.

פרמטרים

יש לנו 4 פרמטרים בתוכנית שלנו. בהינתן אודל השכבה ה-hidden שלנו. בהינתן אודל בתוכנית שלנו. בהינתן $W_1 \in \mathbb{R}^{H \times 784}, b_1 \in \mathbb{R}^H, W_2 \in \mathbb{R}^{10 \times H}, b_2 \in \mathbb{R}^{10}$

כאשר 784 זה גודל הקלט ו-10 זה מספר המחלקות.

בתחילת התוכנית, נרצה ליצור את הפרמטרים האלו ולאתחל אותם בערכים. בקובץ שהמתרגלים התחילת התוכנית, נרצה ליצור את הפרמטרים האלו ומחזיר מטריצה / וקטור מהגודל המתאים העלו הם השתמשו ב-(random.rand(size1, size2) שמחזיר מטריצה – הסתכלו בהערה בהמשך) בטווח (0,1). (לי זה עשה בעיות והייתי צריך לאתחל בצורה אחרת – הסתכלו בהערה בהמשך)

שימו לב! ווקטור בnumpy אפשר לייצג בשתי צורות, לפי ה-shape שלו: אם a.shape מחזיר (1, 784) או (1,784) אז זה בעצם מטריצה עם שורות / עמודות מגודל 1. לעומת זאת אם זה מחזיר (784) אז זה בעצם מטריצה עדיין ndarray שזה מטריצה (784, ית) ואפשר להתייחס אליו (784, 284) אז זה נקרא 1darray (שהוא עדיין משררה או עמודה כנדרש לרוב הדברים. אם לא תספקו לפונקציות את size2 אז הם יצאו בצורה השנייה.

היפר-פרמטרים

אני חושב על 5 היפר-פרמטרים שונים שאפשר לשחק איתם:

- גודל השכבה הנסתרת ה-H שתיארנו קודם הוא גודל שאנחנו קובעים ויכולים לשחק אתו.
 הכלל היחיד בנוגע אליו הוא שהוא יהיה מספר כלשהו בין 784 ל-10 (המטרה היא לאט לאט לצמצם את גודל המידע המודל ילמד איך לעשות את זה בצורה הטובה ביותר (בתקווה))
- קצב הלמידה להזכירכם, כלל העדכון של משקל כלשהו לפי SGD הוא $W_{inew}=W_{iold}-lr*rac{\partial Loss(\Theta)}{\partial W_i}$ זה הלוס עבור הדוגמא הנוכחית בהינתן $W_{inew}=W_{iold}-lr*rac{\partial Loss(\Theta)}{\partial W_i}$ ווא בעצם הגרדיאנט של הלוס לפי W_i ווא כמה "משקל" אנחנו נותנים לגרדיאנט בתיקון הפרמטר. הכלל העיקרי הוא שבד"כ הוא חיובי קטן מ-1. ערכים מומלצים לבדוק הם 0.01, 0.01 ו-0.001 ומשם שחקו כרצונכם.
- מספר האפוקים (epochs) אפוק הוא מעבר על כל הדוגמאות בסט האימון, ואימון עליהן (עדכון משקולות לפי גרדיאנטים של הלוס על הדוגמא). כמה פעמים נעבור על האימון קריטי להצלחת המודל: ישנן דוגמאות שמעבר אחד עליהן לא יצליח ללמד את המודל לסווג אותן נכונה, אבל אם נעשה יותר מדי מעברים על הדאטה נוכל להגיע למצב של overfitting בו המודל למד את סט האימון פרפקט אבל לא ממש עושה חיל בלסווג את סט ה-validation (אפשר לחשוב על זה כאילו המודל למד "לשנן" דוגמאות במקום ללמוד את ההיגיון מאחוריהן)
 - פונקציית אקטיבציה כמו שלמדנו בשיעור, בכדי להכניס אלמנט של אי-ליניאריות למודל
 אנחנו נשתמש בפונקציות לא ליניאריות בתור מעברים בין שכבה לשכבה. ישנן 3 עיקריות
 שאני מכיר: (sigmoid (logistic) שמחזירה תוצאה בין 0 ל-1,

ל-1 ו-ReLU שמחזירה תוצאות אי-שליליות. כל אחת עובדת קצת שונה וטובה לדברים קצת שונים, אבל בשבילנו כולם עובדות די אותו דבר. תנסו כמה שונות ותחליטו מה הכי טוב לכם. שימו לב! יחד עם פונקציית האקטיבציה חשוב להביא את הנגזרת המתאימה לה – זה חלק חשוב מתהליך ה-back propagation שיחשב לנו את הגרדיאנטים.

גודל באטץ' (batch) – עקרונית אפשר לעשות כמה דוגמאות ביחד, לחשב loss משותף ולעדכן פעם ב[גודל באטץ'] דוגמאות (כביכול מקטין את זמן האימון של הרשת – במחיר הדיוק). אני אישית לא כל כך נגעתי בזה וכשנגעתי זה לא הלך לי, אז אני לא אכנס לזה.
 מוזמנים לנסות ③

פסאודו-קוד

אתחול פרמטרים בקטנה...

```
input_size, hidden_size, output_size = 28*28, ..., 10
W1, b1 = initialize_weights(hidden_size, input_size), initialize_weights(hidden_size)
W2, b2 = initialize_weights(output_size, hidden_size), initialize_weights(output_size)
params = [W1,b1,W2,b2]
```

אימון הרשת

נתחיל מסוג של פסאודו-קוד שמתאר את פונקציית האימון ונסביר מה כל חלק אומר.

```
def train(params, epochs, active_func, lr, train_x, train_y, dev_x, dev_y):
    print "epoch", "avg. train loss", "avg. dev loss", "accuracy on dev"
    for i in xrange(epochs): # for each epoch:
        sum loss = 0.0
        shuffle(train x, train y) # shuffle train examples - helps the
model to learn (won't just remember order)
        for x, y in zip(train x, train y):
            out = forward(params, active func, x) # get probabilities
vector as result, where index y is the probability that x is classified as
tag y
            loss = Loss(out, y) # compute loss to see train loss (for
hyper parameters tuning)
            sum loss += loss
            gradients = backprop(params, x, y, active_func) # returns the
gradients for each parameter
            params = update_weights_sgd(params, gradients, lr) # updates
the weights
        dev loss, acc = predict on dev(params, active func, dev x, dev y)
# after each epoch, check accuracy and loss on dev set for hyper parameter
tuning
        print i, sum_loss / train_x.shape[0], dev_loss, "{}%".format(acc *
100)
```

עבור כל אפוק נעבור על כל הדאטה בסט האימון. עבור כל דוגמא נחשב את תוצאת המודל (back propagation), נחשב את הוסגו ואז את הגרדיאנטים של כל משקל (bosk propagation). לאחר (forward pass) מכן נעדכן את המשקלים לפי כלל העדכון של SGD ($W_{inew}=W_{iold}-lr*\frac{\partial Loss(\Theta)}{\partial W_i}$) כמו כן נחשב עבור כל אפוק את הלוס והדיוק על ה-dev_set כדי שנדע מה מצב המודל וכמה ההיפר-פרמטרים שלנו עובדים טוב. נעבור עכשיו על כל חלק ונסביר מה אמור לקרות בו:

הרשת הרשת – Forward Pass

(אני אשתדל להיצמד לשמות שנתנו בתרגול אבל לא מבטיח כלום)

ואת $W_1\in\mathbb{R}^{H\times784}, b_1\in\mathbb{R}^H, W_2\in\mathbb{R}^{10\times H}, b_2\in\mathbb{R}^{10}$ ואת הפרמטרים הבאים: $g\in\{sigmoid,tanh,ReLU\}$ פונקציית האקטיבציה הלא ליניארית

המעבר של קלט ברשת הוא בעצם המשוואות הבאות שיספקו לנו את וקטור ההסתברויות \hat{y} (בהינתן של דוגמא מסוימת):

$$\hat{y} = f(x; W_1, b_1, W_2, b_2) = softmax(z_2) \ (\in \mathbb{R}^{10})$$

$$z_2 = W_2 h + b_2 \ (\in \mathbb{R}^{10})$$

$$h = g(z_1) \ (\in \mathbb{R}^H)$$

$$z_1 = W_1 x + b_1 \ (\in \mathbb{R}^H)$$

לא כזה נורא, נכון? אפשר להסתכל על forward כמעבר על המשוואות האלו מלמעלה למטה, כאשר כל פעם מציבים את מה שיש ומחשבים.

מה זה אומר להפעיל את g על ווקטור? g היא פונקציה שמקבלת x סקלר ומחזירה סקלר אחר, אבל g אנחנו נרצה לקבל וקטור חזרה. פשוט מאוד – הפעלה של g על וקטור תחזיר וקטור שבו על כל איבר הופעל g. למזלנו, בnumpy אפשר לכתוב פונקציות שיעבדו על ווקטורים וגם על ערכים סקלרים – שווה לבדוק את העניין.

שימו לב – ה-back propagation ישתמש בכמה מהערכים שחישבנו ב-forward pass, ולכן אולי אין back propagation אם נחשב אותו בכל מקרה ב-backprop או לחילופין לתת ל-backprop גם את ערכי הביניים הרלוונטיים כדי שלא יחשב שוב. (הערכים הרלוונטיים הם \hat{y} , backprop גם את ערכי הביניים הרלוונטיים לוונטיים מדי שלא יחשב שוב. (הערכים הרלוונטיים הם \hat{y} .

חישוב Loss

ה-loss אומר לנו בעצם כמה מידע "איבדנו" או יותר נכון כמה "רחוקים" אנחנו מלצדוק. המטרה שלנו loss בסופו של דבר היא למזער את ה-loss (ולהעלות את הדיוק, שהם בערך ביחס הפוך אבל לא תמיד). אם נבדוק מה הoss הממוצע ב-train וב-dev כל אפוק נוכל לדעת אם המודל לומד מאפוק לאפוק (וכמה), אם יש overfitting (לוס train יורד אבל dev עולה) ועוד כמה מקרים מעניינים – כלומר סה"כ שווה להשקיע בו.

Back Propagation – הישוב הגרדיאנטים

זה החלק הכי קשוח של התרגיל. אז אם שרדתם עד עכשיו – נשימה עמוקה ונעבור את זה ביחד.

יש לנו 4 פרמטרים, אז אנחנו צריכים 4 גרדיאנטים. אם נסמן את הL-2 loss יש לנו

$$\frac{\partial L}{\partial W_1}$$
, $\frac{\partial L}{\partial b_1}$, $\frac{\partial L}{\partial W_2}$, $\frac{\partial L}{\partial b_2}$

אבל מה זה בעצם אומר גרדיאנט של מטריצה או וקטור? פשוט מאוד, זה מטריצה מגודל זהה שבה כל תא הוא הגרדיאנט של הלוס לפי התא הזה. לדוגמא:

$$gW_1 \coloneqq \frac{\partial L}{\partial W_1} \in \mathbb{R}^{H \times 784}, gW_{1[i,j]} \coloneqq \frac{\partial L}{\partial W_{1[i,j]}}$$

בפועל זה די אומר לגזור כרגיל, ולפעמים צריך לפצל את המטריצות לשורות או ווקטורים לתאים.

 $loss(\hat{y},y) = -\sum_i y_{[i]} \log(\hat{y}_{[i]})$ ושאמרנו שחישוב forward? ושאמרנו את המשוואות היפות מ-forward? הראים למעלה? $\hat{y} = f(x; W_1, b_1, W_2, b_2) = softmax(z_2) (\in \mathbb{R}^{10})$ בשביל לחשב את הגרדיאנטים של כל משקל, אנחנו $z_2 = W_2 h + b_2 \ (\in \mathbb{R}^{10})$ בכלל השרשרת בכדי לחשב אותם בקלות. $h = g(z_1) \ (\in \mathbb{R}^H)$ בכלל השרשרת בכדי לחשב את הגרדיאנטים פה – עזרתי לכמה $z_1 = W_1 x + b_1 \ (\in \mathbb{R}^H)$

(אני לא הולך לחשב את הגרדיאנטים פה – עזרתי לכמה אנשים לחשב אז זה מסתובב, ויש באינטרנט אם אתם לא בטוחים בחישוב. בטוח שתצליחו

נתחיל מ-W2 ומ-b2 – הם יחסית קלים לגזירה, כי הם כמו לגזור שכבה אחת עם סופטמקס (כמו בתרגיל הקודם, אבל שהפעם h הוא הקלט לשכבה). בנוסף בשביל השכבה הבאה, נרצה לחשב את הנגזרת של הלוס לפי הקלט לשכבה הזו (h):

$$\frac{\partial \mathbf{L}}{\partial W_2} = \frac{\partial \mathbf{L}}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial z_2} \frac{\partial z_2}{\partial W_2}, \qquad \frac{\partial \mathbf{L}}{\partial b_2} = \frac{\partial \mathbf{L}}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial z_2} \frac{\partial z_2}{\partial b_2}, \qquad \frac{\partial \mathbf{L}}{\partial h} = \frac{\partial \mathbf{L}}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial z_2} \frac{\partial z_2}{\partial h}$$

ואלו הגרדיאנטים של כל השכבה השנייה. נשמור הצד את הגרדיאנטים של W2 ו-b2, ונמשיך לחשב b2. את של W1 ו-b2.

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial h} \frac{\partial h}{\partial z_1} \frac{\partial z_1}{\partial W_1}, \qquad \frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial h} \frac{\partial h}{\partial z_1} \frac{\partial z_1}{\partial b_1}$$

נשים לב ש- $\frac{\partial h}{\partial z_1}=g'(z_1)$ כלומר לנגזרת של הפונקציה g שבחרנו (משתנה לפי בחירת g). נשים לב ש- לב $z_1=g'(z_1)$ בחירת של backprop(params, x, y, active_func) שהפונקציה (active_func בכדי לחשב את הגרדיאנטים.

(שוב, לחשב את הנוסחאות לנגזרות זה אחד החלקים הקשים בתרגיל. אל תתייאשו, וקחו הרבה השראה מהמצגות – הן עושות אותו תהליך פשוט על רשת טיפה שונה)

?למה בעצם מחשבים גרדיאנטים

אז למי שלא הבין למה הוא גוזר פונקציות אחרי אינפי 2 (ברוך שפטרנו) - זה בשביל לצמצם את המי למי שלא הבין למה הוא גוזר פונקציות אחרי אינפי 2 (ברוך שפטרנו) - זה בשביל לצמצם את המוא למי למי לצמצם את

בהרצאות הראשונות הגדרנו את פונקציית הלוס שבהינתן תוצאת הרשת עבור דוגמה x כלשהי ואת התיוג האמתי שלה y ומחזירה כמה מידע "איבדנו" או כמה "רחוקים" אנחנו מלצדוק.

עכשיו המטרה בלמידה היא לצמצם את הלוס על כל דוגמא, ע"י שנמצא את הפרמטרים (W-ים ו-b-טים) שיתנו את הלוס הכי נמוך על כל הדוגמאות - בעצם יש לנו פונקציה (loss) שאנחנו רוצים למצוא (את נקודת המינימום עבורה (כאשר ה"נקודה" היא ערכי ה-W-ים ו-b-ט).

עכשיו אם היית נותן לכם פונקציה לא ידועה ומבקש למצוא לה את המינימום, הייתם יכולים לנסות לחשב ערך בנקודה ואז לחשב את הנגזרת בנקודה - אם נקבל נגזרת חיובית נדע שאנחנו בעליה ונרצה ללכת בכיוון השני, ואם שלילית נדע שכדאי ללכת בכיוון בשביל למצוא מינימום.

זה פישוט עצום למה שקורה בלמידה - האלגו' SGD דואג לצמצם לאט לאט את המשקלים לכיוון המינימום, באמצעות הIoss שיצא ("כמה טעינו") והגרדיאנטים של הIoss לפי כל פרמטר ("באיזה כיוון צריך לתקן"). (מי שלמד בינה אולי זוכר שגם אז דיברנו על "כמה טעינו" ו"באיזה כיוון") עד שנתכנס לערכים טובים יחסית.

עדכון המשקלים

פשוט מאוד לעדכן את המשקלים לפי SGD (זהו נשבע פעם אוד לעדכן את המשקלים לפי און אורונה שאתם רואים בקובץ את הנוסחה הזו אחרונה שאתם רואים בקובץ את הנוסחה ה

```
dev set-על ה-loss חישוב דיוק ו
```

סיימנו אפוק אחד של למידה על הדוגמאות! כל הכבוד לנו!

אבל לפני שמתחילים לפתוח שמפניות בואו נוודא שזה באמת עובד לנו – ובדיוק בשביל זה הכינינו מראש Validation Set.

הרעיון פשוט מאוד: נעבור על כל הדוגמאות, ועבור כל אחת נחשב את התוצאה מה-forward ואת ה-loss, נספור כמה פעמים צדקנו ונחזיר את ממוצע הלוס ואת הדיוק.

נשים לב שאין צורך ב-backprop ולא בעדכון משקולות – למעשה זה אפילו **אסור!** המטרה שלנו היא בעצם לא ללמוד את הדוגמאות ב-dev, אלא לבדוק מה מצבנו. זה הכל.

:פסאודו-קוד קטנטנן

```
def predict_on_dev(params, active_func, dev_x, dev_y):
    sum_loss = good = 0.0  # good counts how many times we were correct
    for x, y in zip(dev_x, dev_y):
        out = forward(params, active_func, x)  # get probabilities
vector as result, where index y is the probability that x is classified
as tag y
        loss = Loss(out, y)  # compute loss to see train loss (for
hyper parameters tuning)
        sum_loss += loss
        if out.argmax() == y:  # model was correct
            good += 1
        acc = good / dev_x.shape[0]  # how many times we were correct / #
of examples
    avg_loss = sum_loss / dev_x.shape[0]  # avg. loss
    return avg_loss, acc
```

ובכן, כמעט. כלומר, סיימנו לכתוב את רוב הקוד, אבל עדיין יש דברים לעשות – למצוא את ההיפר פרמטרים שיהפכו את הרשת שלנו להיות הרשת הכי טובה שיש.

איך עושים את זה? פשוט מאוד, מנסים שינויים ובודקים מה עובד יותר טוב ומה פחות. הדיוק הוא הפרמטר החשוב כי בסופו של דבר מעניין אותנו להיות כמה שיותר צודקים – אבל הלוס חשוב כמעט באותה מידה כי כשהלוס יורד זה אומר שהרשת לומדת.

הרעיון הוא שתמשיכו לנסות ולהריץ עד שמגיעים לתוצאה שאוהבים – כי השלב הבא הוא המבחן האמתי.

Predicting on the Test Set

כעת הגיע הרגע הגורלי, בו אנחנו נוציא מהרשת פלטים שעל פיהם יקבע הציון שלנו. זה מאוד פשוט בפועל – נעבור על כל דוגמא ב-test_x, נעשה עבורה forward וניקח את ה-(אגב למה לוקחים argmax)? זה בעצם הקלאס שהמודל חושבת שיש את ההסתברות הגבוהה ביותר שהוא הנכון) ונכתוב אותו לשורה המתאימה בקובץ.

יאללה פסאודו-קוד אחרון לכיף:

```
f = open("test.pred", "w")
for x in test_x:
    out = forward(params, active_func, x)
    f.write(str(out.argmax()))
f.close()
```

בעיות, פתרונות וטיפים

אני אנסה לתאר פה כל מיני בעיות שאני נתקלתי בהן ופתרונות נפוצים שאולי הם לא טריוויאליים. אם למישהו קרה משהו מעניין הוא מוזמן להגיד לי ואני אוסיף.

"פיצוץ" של תוצאות השכבה הראשונה – לי קרה הרבה בהתחלה ש-z1 היו מספרים ממש גדולים, ואז sigmoid או למיד 1 על כולם, וזה דפק את הפרדיקציה ואת הגרדיאנטים.

הפתרון שמצאתי הוא דווקא שלושה:

1) נרמול הקלטים: (בהצעתו של תמיר המלך) נחלק את כל הקלטים ב-255.0, ככה שבמקום קלטים בטווח 0-255.0 נעבור לקלטים בטווח [0,1].

 $(W \sim U(-0.08, 0.08))$ אתחול המשקלים בצורה יוניפורמיות בטווח בין 0.08 ל-2. ל-20.08 (2

3) אתחול על פי השיטה של Glorot: בקצרה, בהינתן שאנחנו רוצה לאתחל משקל בגדלים

 $W \sim U(-arepsilon,arepsilon)$ ידגם יוניפורמית: W ידגם ואז המשקל פ $arepsilon = \sqrt{rac{6}{n+m}}$ או א n x m

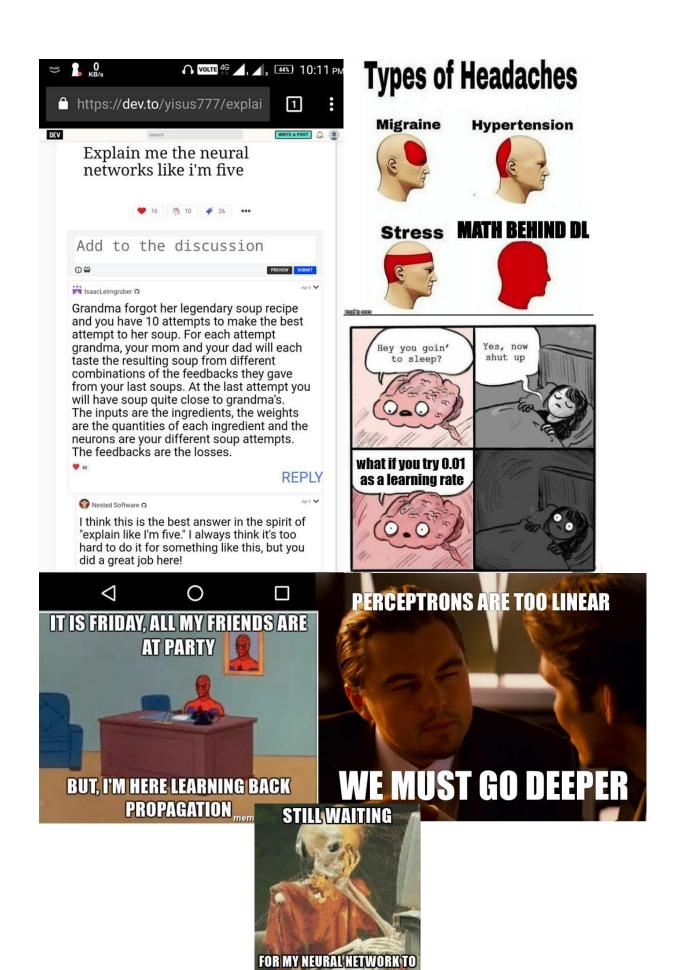
(random.uniform בשם numpy יש פונקציית)

▶ Loss אינסופי – זה קורה בד"כ בגלל שהמודל "ממש בטוח" בעצמו אז התוצאה מהסופטמקס היא ווקטור של אפסים עם אחד במקום שהוא חושב. אז אם הוא טועה, הלוס שהוא מקבל על להגיד שיש %0 סיכוי שהתיוג הוא התשובה הנכונה הוא גדול מאוד – ולמעשה אינסופי (כי (log(0)...)

פתרון: כרגע אין. אשמח לשמוע הצעות. בעיקר כרגע הפסקתי להשתמש ב-ReLU כי הוא עשה לי את הבעיות האלו.

רשת איטית – הרשת לא פשוטה למימוש לכן סביר להניח שלא נממש אותה בצורה האופטימלית וחבל. יש הרבה פונקציות של numpy שיכולות לעזור לנו בביצועים (למשל dot, outer) ושימוש בפעולות שעובדות על כל המטריצה).

באופן כללי קחו לכם אתגר: בלי להשתמש בלולאות בכלל, חוץ מה-4 שצריך בשביל לעבור על הדאטה (מעבר על כל אפוק, מעבר על כל דוגמא בסט האימון, מעבר על כל סט הvalidation ומעבר על כל סט ה-test). ממליץ להיעזר ב tutorial שיש בהקדמה.



TRAIN