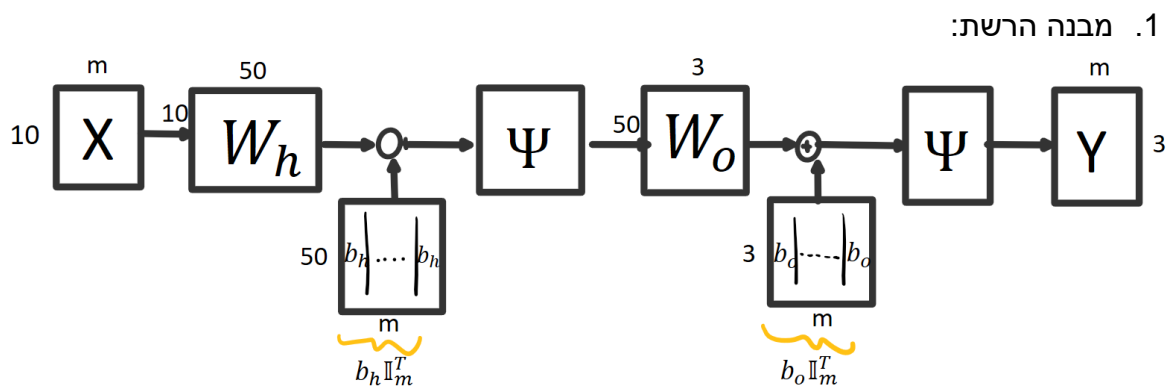


# תרגיל בית 1

מגשים:

שם	ת.ז.
גבריאל חביב	
איטן חצרוני	302383856

## חלק תיאורטי

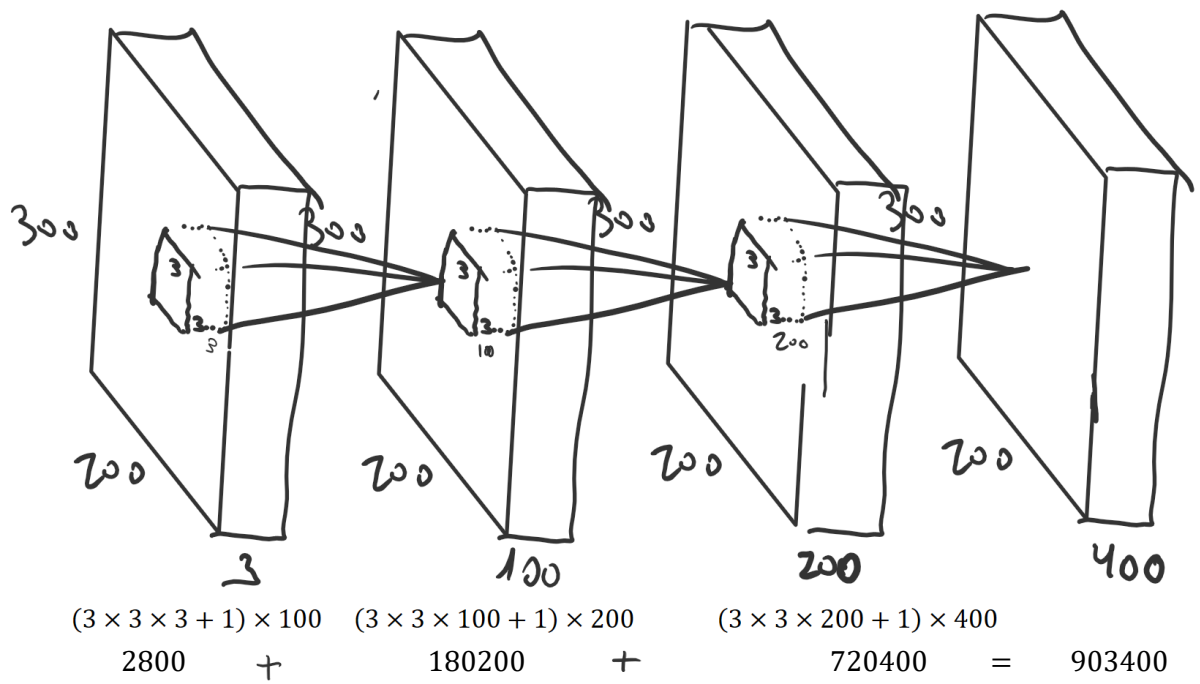


- המימדים של  $X$  הם  $10 \times m$ . זהו גודל וקטור הכניסה, ו- $m$  גודל ה-batch.
- המימדים של  $W_h$  הם  $10 \times 50$ . המימד של וקטור ה-bias,  $b_h$ , הוא  $50 \times 1$ . יש להוסיף את אותו ה-bias לכל התמונות, לכן למעשה צריך להוסיף מטריצה בגודל  $50 \times m$ , כאשר היא מורכבת מעמודות של הוקטור  $b_h$ . ניתן לכתוב זאת בתור  $b_o \mathbb{I}_m^T$ , כאשר  $\mathbb{I}_m$  הוא וקטור אחדות בגודל  $m$ .
- המימדים של  $W_o$  הם  $50 \times 3$  ושל וקטור ה- $b_o$  הם  $3 \times 1$ . גם כאן יש להוסיף מטריצה בגודל  $3 \times m$  שעמודותיה בנויות מהוקטור  $b_o$ . ניתן לכתוב זאת בתור  $b_o \mathbb{I}_m^T$ .
- הצורה של  $Y$  היא מטריצה בגודל  $3 \times m$ .

$$e. Y = \Psi(W_o^T \Psi(W_h^T X + b_h \mathbb{I}_m^T) + b_o \mathbb{I}_m^T)$$

כאשר  $\Psi$  היא פעולת ReLU איבר איבר על המטריצה

2. מבנה הרשת:



סה"כ ישנם 903400 פרמטרים.

$$3. a. \frac{\partial f}{\partial r} = \sum_{i=1}^m \frac{\partial f}{\partial y_i} \cdot \frac{\partial y_i}{\partial r} = \sum_{i=1}^m \frac{\partial f}{\partial y_i} \cdot \frac{\partial}{\partial r} (r \hat{x}_i + \beta) = \sum_{i=1}^m \frac{\partial f}{\partial y_i} \cdot \hat{x}_i$$

$$b. \frac{\partial f}{\partial \beta} = \sum_{i=1}^m \frac{\partial f}{\partial y_i} \cdot \frac{\partial y_i}{\partial \beta} = \sum_{i=1}^m \frac{\partial f}{\partial y_i} \cdot \frac{\partial}{\partial \beta} (r \hat{x}_i + \beta) = \sum_{i=1}^m \frac{\partial f}{\partial y_i}$$

$$c. \frac{\partial f}{\partial \hat{x}_i} = \frac{\partial f}{\partial y_i} \cdot \frac{\partial y_i}{\partial \hat{x}_i} = \frac{\partial f}{\partial y_i} \cdot \frac{\partial}{\partial \hat{x}_i} (r \hat{x}_i + \beta) = \frac{\partial f}{\partial y_i} \cdot r$$

$$d. \frac{\partial f}{\partial r^2} = \sum_{i=1}^m \frac{\partial f}{\partial y_i} \cdot \frac{\partial y_i}{\partial r^2} = \sum_{i=1}^m \frac{\partial f}{\partial y_i} \cdot \frac{\partial y_i}{\partial \hat{x}_i} \cdot \frac{\partial \hat{x}_i}{\partial r^2} = \sum_{i=1}^m \frac{\partial f}{\partial y_i} \cdot \frac{\partial}{\partial \hat{x}_i} (r \hat{x}_i + \beta) \cdot \frac{\partial}{\partial r^2} \left( \frac{x_i - \mu_B}{\sqrt{r_B^2 + \epsilon}} \right)$$

$$= -\frac{1}{2} r \sum_{i=1}^m \frac{\partial f}{\partial y_i} \cdot \frac{x_i - \mu_B}{(r_B^2 + \epsilon)^{3/2}} = -\frac{r}{2(r_B^2 + \epsilon)^{3/2}} \sum_{i=1}^m \frac{\partial f}{\partial y_i} (x_i - \mu_B)$$

$$e. \frac{\partial f}{\partial \mu_B} = \sum_{i=1}^m \frac{\partial f}{\partial y_i} \cdot \frac{\partial y_i}{\partial \mu_B} = \sum_{i=1}^m \frac{\partial f}{\partial y_i} \cdot \frac{\partial y_i}{\partial \hat{x}_i} \cdot \frac{\partial \hat{x}_i}{\partial \mu_B} = \sum_{i=1}^m \frac{\partial f}{\partial y_i} \cdot r \cdot \frac{\partial}{\partial \mu_B} \left( \frac{x_i - \mu_B}{\sqrt{r_B^2 + \epsilon}} \right) =$$

$$r \sum_{i=1}^m \frac{\partial f}{\partial y_i} \left( \frac{-1}{\sqrt{r_B^2 + \epsilon}} - \frac{x_i - \mu_B}{2(r_B^2 + \epsilon)^{3/2}} \cdot \frac{\partial r_B^2}{\partial \mu_B} \right) = r \sum_{i=1}^m \frac{\partial f}{\partial y_i} \left( \frac{-1}{\sqrt{r_B^2 + \epsilon}} + \frac{x_i - \mu_B}{2(r_B^2 + \epsilon)^{3/2}} \cdot \frac{-2}{m} \sum_{j=1}^m (x_j - \mu_B) \right)$$

$m \cdot \frac{1}{m} \sum x_i - m \mu_B = 0$

$$= -\frac{r}{\sqrt{r_B^2 + \epsilon}} \sum_{i=1}^m \frac{\partial f}{\partial y_i}$$

$$f. \frac{\partial f}{\partial x_i} = \frac{\partial f}{\partial \hat{x}_i} \cdot \frac{\partial \hat{x}_i}{\partial x_i} + \frac{\partial f}{\partial \mu_B} \cdot \frac{\partial \mu_B}{\partial x_i} + \frac{\partial f}{\partial r_B^2} \cdot \frac{\partial r_B^2}{\partial x_i} = r \frac{\partial f}{\partial y_i} \cdot \frac{1}{\sqrt{r_B^2 + \epsilon}} - \frac{r}{\sqrt{r_B^2 + \epsilon}} \sum_{j=1}^m \frac{\partial f}{\partial y_j} \cdot \frac{1}{m}$$

$$- \frac{r}{2(r_B^2 + \epsilon)^{3/2}} \sum_{j=1}^m \frac{\partial f}{\partial y_j} (x_j - \mu_B) \cdot \frac{\partial}{\partial x_i} (x_j - \mu_B)$$

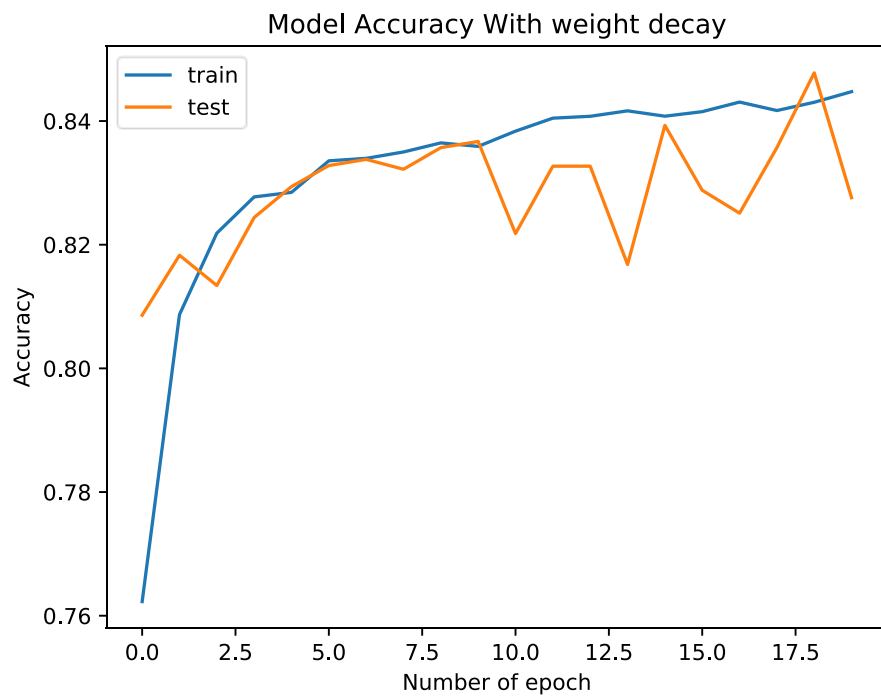
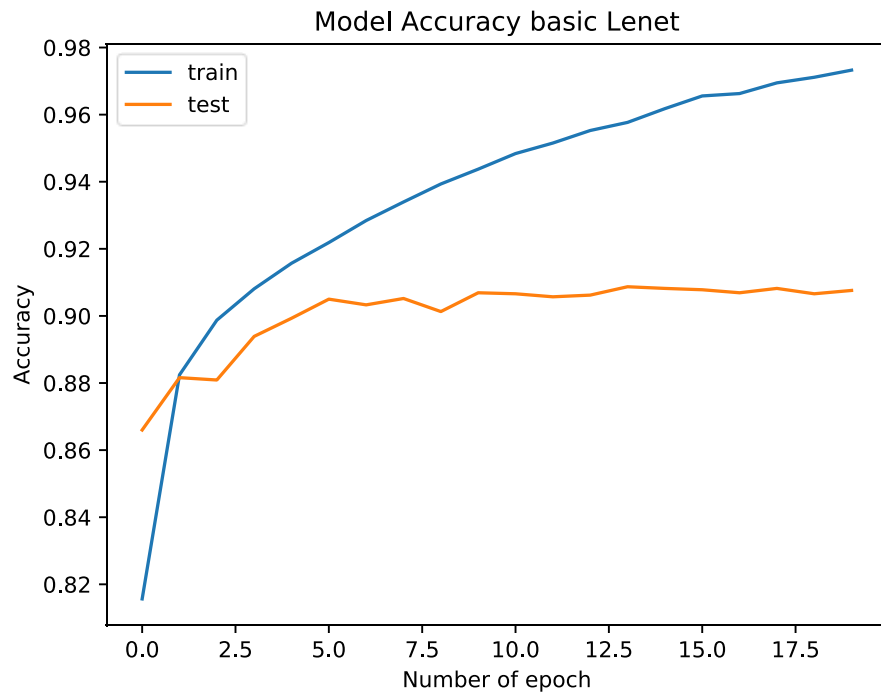
$$= \frac{r}{m \sqrt{r_B^2 + \epsilon}} \left( m \frac{\partial f}{\partial y_i} - \sum_{j=1}^m \frac{\partial f}{\partial y_j} - \frac{1}{\sqrt{r_B^2 + \epsilon}} (x_i - \mu_B) \sum_{j=1}^m \frac{\partial f}{\partial y_j} \frac{(x_j - \mu_B)}{\sqrt{r_B^2 + \epsilon}} \right)$$

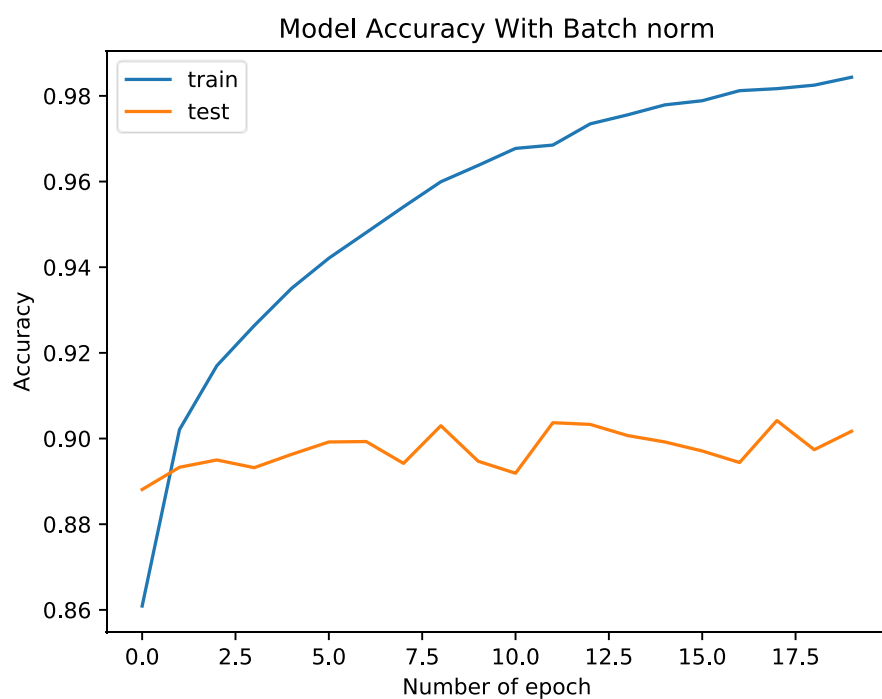
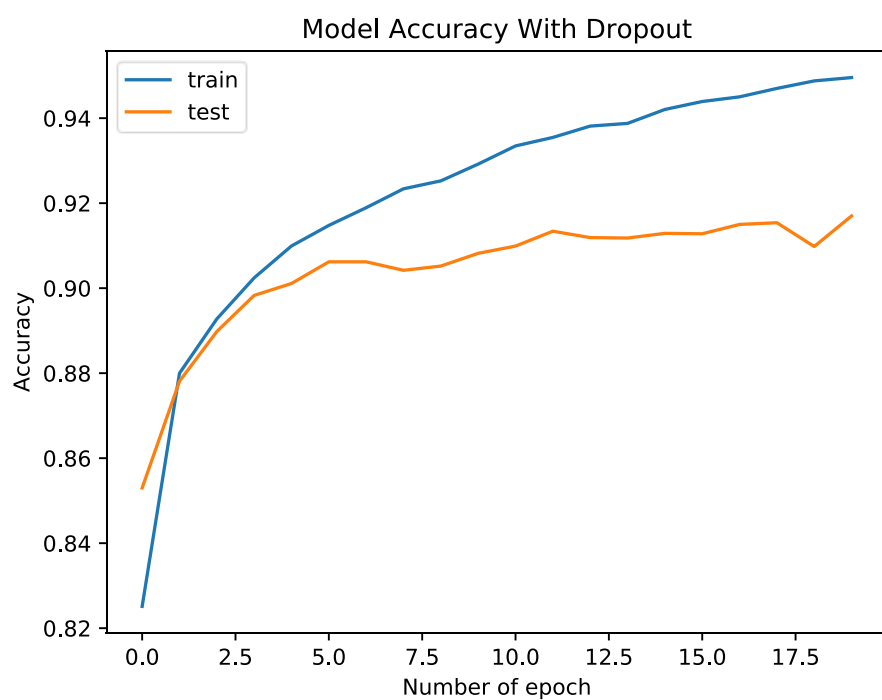
$\hat{x}_i$        $\hat{x}_j$

$$= \frac{r}{m \sqrt{r_B^2 + \epsilon}} \left( m \frac{\partial f}{\partial y_i} - \sum_{j=1}^m \frac{\partial f}{\partial y_j} - \hat{x}_i \sum_{j=1}^m \frac{\partial f}{\partial y_j} \hat{x}_j \right)$$

## חלק מעשי

הגרפים המבוקשים:





סיכום ביצועים:

Accuracy	טכניקה
0.91	Normal
0.92	Dropout
0.83	Weight Decay
0.90	Batch Normalization

## הערות:

1. ניתן לראות הביצועים של רוב הטכניקות דומים.
2. הביצועים של טכניקת weight decay נמוכים באופן יחסי, אך עם שגיאת ההכללה הנמוכה ביותר. הסיבה לכך היא שאמנם הרשת לא מגיעה למצב של over-fitting, אך הרגולריזציה מונעת ממנה להיות מורכבת מספיק בכדי לסווג את המידע. ניתן ככל הנראה בעזרת אופטימיזציה של איבר הרגולריזציה להגיע ל-trade-off מוצלח יותר בין ביצועים לשגיאת הכללה.
3. ניתן לראות שבעזרת batch normalization הרשת התכנסה לערך סופי בצורה מהירה, ומשם גדלה שגיאת ההכללה (וה-accuracy כמעט ולא השתנה).