

מערכות לומדות תשפ"ג - תרגיל 1

התפלגויות, קונבולוציה וחידה

בתרגיל זה תממשו בקוד ותבחנו כמה נושאים המתוארים בשיעור 2 – מבוא להסתברות, וכן תרחיבו את היכרותכם עם אופרטור הקונבולוציה ושימושו בעיבוד אותות, ולסיום תענו על חידת הקלסיפיקציה.

סעיפים להגשה במסמך הפתרון \ קבצי קוד מסומנים בצבע חום. סעיפים לברור עצמי מסומנים בירוק.

נושא 1 - פעולת הקונבולוציה, פעולות קרובות ומספר שימושים

בשיעור 2 תארנו משתנה מקרי בדיד ורציף ואת ההתפלגויות שלהם. בתנאים מסוימים ההתפלגות של סכום משתנים מקריים בלתי תלויים מתקרבת להתפלגות נורמלית (גאוסיאנית).

בשאלה זו תחקרו את פעולת הקונבולוציה בהגדרתה המתמטית עבור המקרה הבדיד, ותשתמשו בה ככלי עזר בביצוע פעולות של עיבוד אותות, זאת כהכנה לחילוץ מאפיינים בהמשך הקורס.

ההגדרת קונבולוציה למקרה הבדיד:

$$\begin{aligned}(f * g)[n] &\stackrel{\text{def}}{=} \sum_{m=-\infty}^{\infty} f[m] g[n - m] \\ &= \sum_{m=-\infty}^{\infty} f[n - m] g[m].\end{aligned}$$

שימו לב שזו הגדרה סימטרית ל f ול g .

לפעמים נרצה להשתמש בפונקציות שתחום הערכים שלהן מוגבל - finite support, ואף אינו באותו אורך. זה נפוץ כאשר g מהווה את גרעין הקונבולוציה ותחום הערכים שלה קטן בצורה משמעותית משל f .

עבור מקרה זה של תחום ערכים מוגבל עבור g

$$\{-M, -M + 1, \dots, M - 1, M\}$$

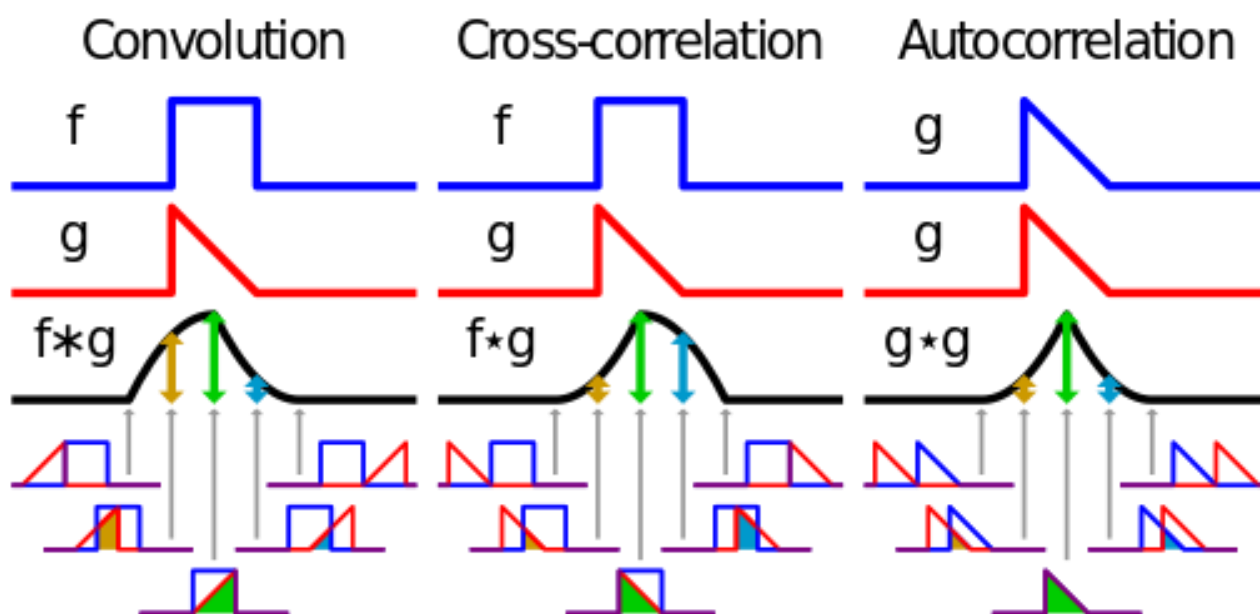
ההגדרה היא:

$$(f * g)[n] = \sum_{m=-M}^M f[n - m] g[m]$$

בררו לעצמכם:

- כמה איברים יש בתוצאה עבור ערך n מסוים?
בקוד המבצע קונבולוציה, כפי שממומש בחבילות תוכנה שונות, מחשבים מספר ערכים על ידי ביצוע הפעולה עבור כל ערכי n האפשריים בהינתן האותות.
- עבור n מסוים, בהנחה שעוברים על כל האיברים של האות g בסכום, כמה איברים יש לאות g ?
- רשמו לעצמכם את ההגדרה כאשר m רץ מ 0 ועד לאיבר האחרון ב g . זו ההגדרה המתאימה למעבר על פני כל איברי מערך g .
- בביצוע הסכום מתקדמים על פני האות g עם העלייה בערך אינדקס m . מה כיוון התנועה על פני האות f ?

לעיתים משתמשים בפעולות הקשורות לקונבולוציה: קרוס-קורלציה ואוטו-קורלציה. הנה המחשה לפעולתן (מתוך וויקיפדיה)



בשיעור הזכרנו משפט הטוען כי התפלגות של סכום שניים (או יותר) משתנים מקריים **בלתי תלויים** היא הקונבולוציה של ההתפלגויות:

https://en.wikipedia.org/wiki/Convolution_of_probability_distributions

הנה רשימה של התפלגויות ידועות וההתפלגות של הסכום שלהן:

https://en.wikipedia.org/wiki/List_of_convolution_of_probability_distributions

בנוסף הזכרנו את משפט הגבול המרכזי בהקשר של סכום משתנים מקריים - סכום של משתנים מקריים בלתי תלויים נוטה להתפלג נורמלית ככל שמספר המשתנים המקריים גדול יותר (גם אם המשתנים המקריים עצמם אינם מתפלגים נורמלית!):

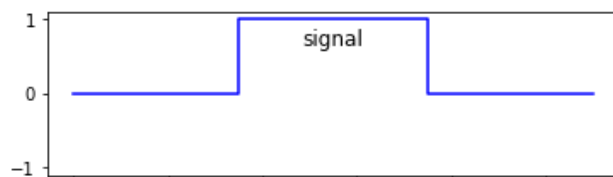
https://en.wikipedia.org/wiki/Central_limit_theorem

קונבולוציה מהווה כלי עזר חישובי ליצירת ההתפלגות של הסכומים.

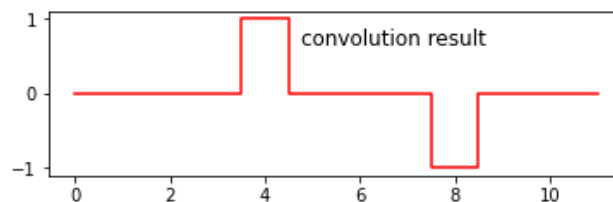
כעת נשתמש בקונבולוציה למשימה אחרת – עיבוד אותות.

שאלה 1

א. המערך החד ממדי signal מכיל 12 איברים: $[0,0,0,0,1,1,1,1,0,0,0,0]$



לפניכם גרף של הערכים **בכחול**. זו הצגה של האות שהוא אחד הקלטים לפעולת הקונבולוציה, הקלט האחר הוא הגרעין kernel המכיל שני איברים $[a,b]$,



וגרף של תוצאת הקונבולוציה **באדום**.

להגשה:

א1. ענו במסמך הפתרון: מה ערכי a, b של ה kernel הנדרשים ליצירת תוצאה זו?

א2. קובץ קוד בשם ex1_a_1.py המחשב קונבולוציה של ה signal ו ה kernel ומייצר שלושה גרפים זה מעל זה: הצגת ה signal, ה kernel והתוצאה.

לצורך הצגת ה kernel השתמשו במערך $[0,a,b,0]$.

על הגרפים להראות בדיוק על פי הדוגמה שכאן, כאשר גרף ה kernel יהיה אמצעי והעקומה תוצג בירוק. שימו לב שלכל הגרפים אותו טווח ערכים לצירים.

א3. העתיקו את הגרפים שהקוד שלכם מייצר למסמך הפתרון.

עזרה:

השתמשו ב np.convolve לביצוע קונבולוציה. הסתכלו בהסבר של numpy על הפונקציה הזו והפרמטרים שלה.
השתמשו ב matplotlib לצרכי התצוגה. הסתכלו בתיעוד של subplots, set_title, step.

ב. טענה:

הפעלת הקונבולוציה עם ה kernel של השאלה הקודמת על אות חד ממדי היא ביצוע של גזירה דיסקרטית (discrete derivative) של האות, כלומר ביצוע קרוב דיסקרטי לפעולת הנגזרת.

להגשה:

1. הסבירו את הטענה בעזרת הגדרת הנגזרת. התייחסו במפורש ל h ומשמעותו בקרוב דיסקרטי של הנגזרת.

$$f'(a) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h}$$

2. הסבירו:

מדוע יש צורך בקרוב דיסקרטי של נגזרת? האם לא עדיף להשתמש ישירות בהגדרה המתמטית עם הגבול שהוצגה כאן?

ג. מצאו גרעין שהפעלתו תהווה קרוב דיסקרטי לנגזרת שנייה. עשו זאת על פי השלבים הבאים:

A.

אם f מייצגת את הגרעין לגזירה (משאלה 1א) ו d מייצגת את הנתונים אזי נגזרת ראשונה בעזרת קונבולוציה תיכתב כך

$$f*d$$

B.

נגזרת שנייה היא הפעלת פעולת הנגזרת פעמיים:

$$f'' = (f')'$$

על כן בעזרת קונבולוציה נרשום כך:

$$f*(f*d)$$

.C

פעולת הקונבולוציה היא פעולה אסוציאטיבית.
על כן גזירה פעמיים נרשום כך

$$f*(f*d) = (f*f)*d$$

נשים לב ש $f*f$ היא הפעלת קונבולוציה של הגרעין עם עצמו.

נכנה את התוצאה של פעולה זו בשם $f2$ וזה יהיה הגרעין לביצוע פעולת נגזרת שנייה

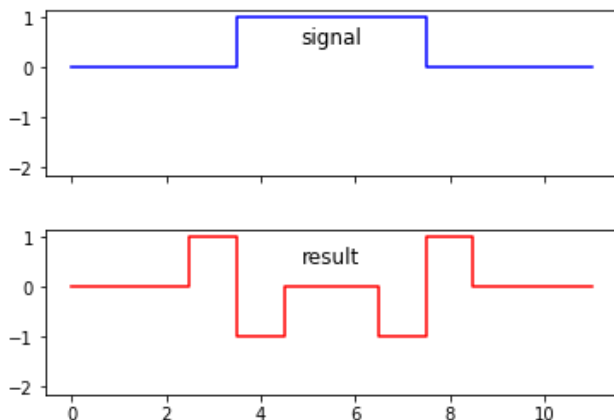
$$f*(f*d) = (f*f)*d = f2*d$$

ג1. להגשה: מצאו את הערכים של $f2$ ורשמו את התוצאה.
(אפשר להשתמש בפונקציית הקונבולוציה של `numpy.convolve` כדי למצוא את הערכים, או לבצע קונבולוציה או קרוס-קורלציה ידנית).

הערה: הגרעינים אותם אנו מציגים כאן אינם יחידים, ויש עוד קרובים דיסקרטיים לנגזרות. ראו למשל כאן:

<https://towardsdatascience.com/image-derivative-8a07a4118550>

כתבו קוד המבצע נגזרת שניה על ה `signal` משאלה 1א ומציג את התוצאה על ידי שלושה גרפים בדומה לשאלה 2א1



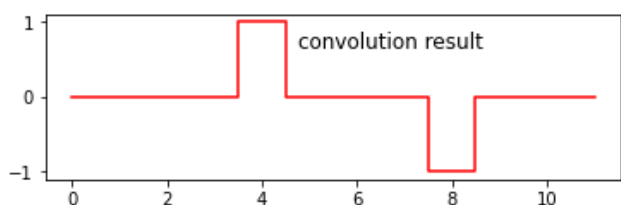
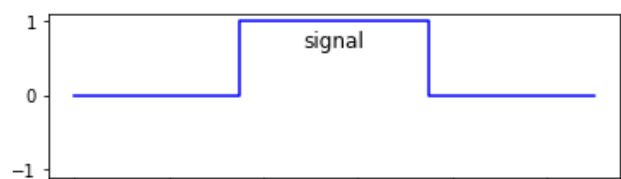
הנה התוצאה של הקונבולוציה:

ג2. להגשה: קובץ קוד בשם `ex1_c_2.py` המחשב קונבולוציה של ה `signal` ו ה `f2 kernel` של נגזרת שניה ומייצר שלושה גרפים זה מעל זה: הצגת ה `signal`, ה `kernel` והתוצאה.
לצורך הצגת ה `kernel` השתמשו במערך `[0, f2, 0]`.

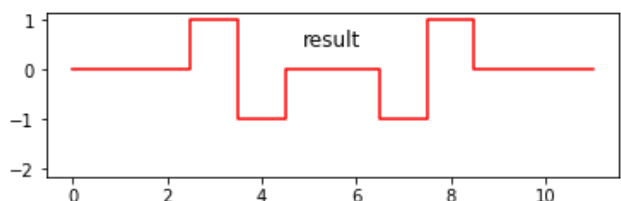
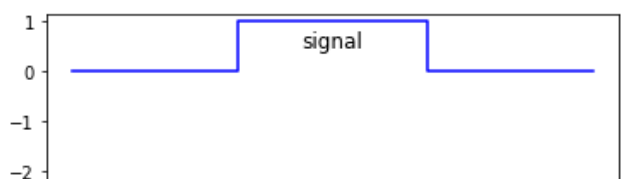
על הגרפים להראות בדיוק על פי הדוגמה שכאן, כאשר גרף ה kernel יהיה אמצעי והעקומה תוצג בירוק. שימו לב שלכל הגרפים אותו טווח ערכים לצירים.

ג3. העתיקו את הגרפים שהקוד שלכם מייצר למסמך הפתרון.

הסבר לפעולות שהוצגו עד כה:



הנגזרת הראשונה מדגישה קצוות (edges), כלומר אזורים בהם יש מעבר בין ערכים. למשל כאן, המעברים הם בשני קצות "התיבה" הכחולה. המעבר השמאלי, בין ערך נמוך לערך גבוה מניב בתוצאה ערך חיובי. המעבר הימני, בין ערך גבוה לערך נמוך מניב בתוצאה ערך שלילי. הסימן של התוצאה בפני עצמו אינו כה חשוב (הרי נקבל סימן ההפוך אם נשתמש בסדר ערכים הפוך בגרעין). אבל תמיד נקבל סימן תוצאה הפוך בתוצאה של שני קצות התיבה.



הנגזרת השנייה גם היא מדגישה קצוות. אך שימו לב שכעת, לכל קצה באות המקור מתקבל שינוי כפול בתוצאה. גם כאן סימן הערכים של התוצאה בשני הקצוות הפוך. התוצאה המתקבלת על ידי גזירה שנייה בעזרת קונבולוציה עם f_2 הייתה מקבלת גם אם היינו גוזרים שוב (קונבולוציה עם f) את התוצאה של הנגזרת הראשונה (הגרף האדום שמעל).

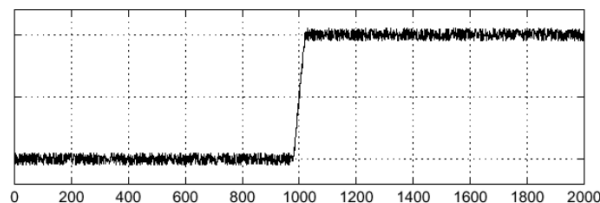
אם המשימה היא זיהוי קצוות (שינויים), מה היתרון בשימוש בנגזרת שנייה בעיבוד אותות ועיבוד תמונה?

כדי להסביר נזכור שבמציאות, רוב האותות כוללים גם רעש ואינם כה ברורים כמו התיבה הכחולה. נדגים זאת בעזרת

[https://www.cs.cmu.edu/~16385/s17/Slides/4.0 Image Gradients and Gradient Filtering.pdf](https://www.cs.cmu.edu/~16385/s17/Slides/4.0%20Image%20Gradients%20and%20Gradient%20Filtering.pdf)

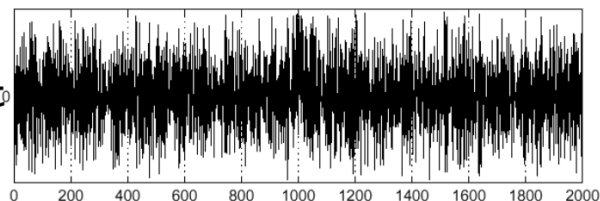
הנה אות חד ממדי רועש. ננסה לזהות בו את הedge, השינוי הגדול בגובה סביב $x=1000$:

Intensity plot



לצורך זה נגזור בעזרת קונבולוציה עם פילטר גזירה:

Derivative plot



מה קרה?

הבעיה היא קנה מידה. השינוי הגדול בגובה מתבטא בקנה מידה לא קטן בציר ה X , הוא נפרש על פני כ 10 ערכי x סמוכים, ואילו הרעשים מתבטאים בקנה מידה קטן בהרבה בציר ה X .

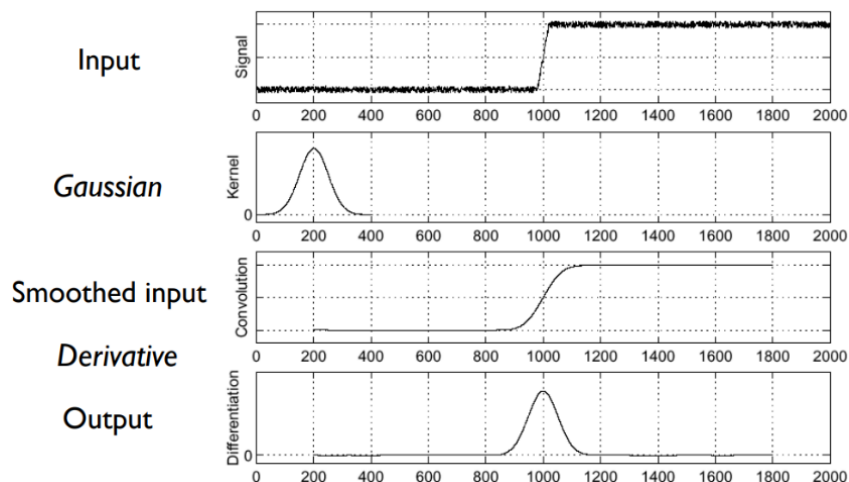
פילטר הגזירה שהשתמשנו בו כאן מתאים לקנה מידה קטן: הוא מורכב מ 2 או 3 ערכים בלבד ולכן הקונבולוציה בעזרתו מדגישה תופעות הפרושות על קנה מידה קטן בציר ה X . במילים אחרות, הגזירה מדגישה את הרעש, ולא מוצאת את השינוי באות.

מה הפתרון?

אפשר לבנות פילטר גזירה ארוך יותר (המכיל יותר איברים), ובכך לנסות להגיע לקנה המידה הנכון. דרך נוחה יותר היא להפריד רעיונית בין שתי המשימות: 1. הקטנת הרעש בקנה המידה המתאים. 2. ביצוע גזירה לזיהוי השינויים הנותרים.

העלמת הרעש מתבצעת בעזרת פילטר החלקה. אותו נבנה בקנה המידה המתאים: כך שיעילים ככל האפשר את הרעש אך לא ישנה מאוד את התופעה בקנה המידה הגדול יותר. פילטר מועיל ומקובל להחלקה הוא גאוסיאן (נקרא לו G).

הנה רצף הפעולות:



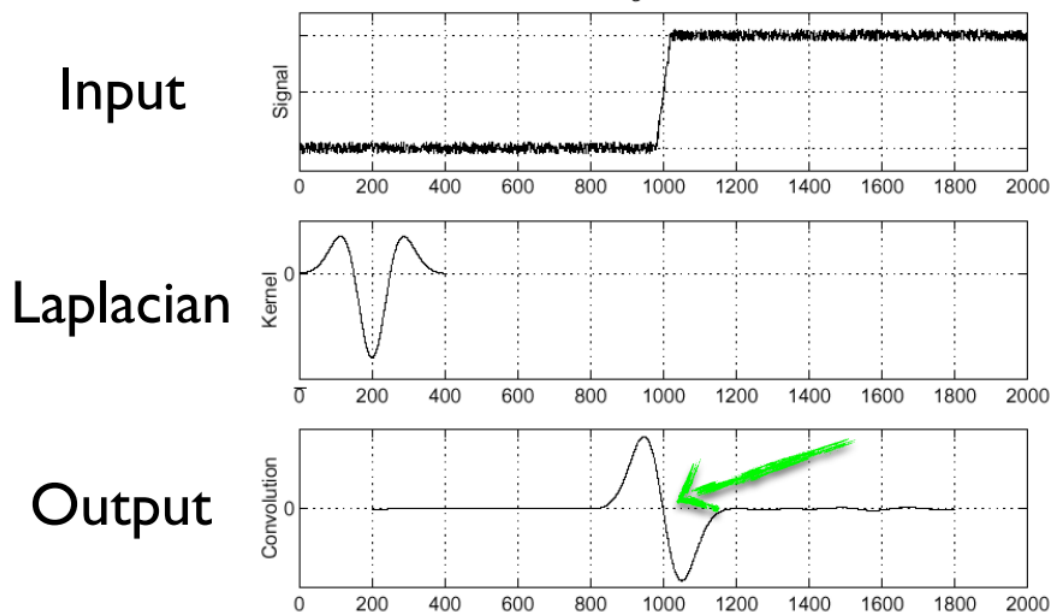
כעת, כדי לזהות את מיקום ה edge באות המקורי, נחפש את ערך המקסימום בתוצאה.

בזכות האסוציאטיביות של הקונבולוציה, במקום לבצע שתי פעולות קונבולוציה על האות המקורי, אפשר כמובן לגזור את פילטר הגאוסיאן G ולקבל פילטר חדש dG (איך הוא נראה?), ואז לבצע קונבולוציה על האות המקורי בעזרת dG . התוצאה הסופית תהיה זהה, אך עבור אות מקורי גדול הפעולה תהיה מהירה יותר (מדוע?).

נחזור לשימוש בנגזרת שנייה:

אפשר להחליק את האות בעזרת גאוסיאן G ואז לגזור אותו נגזרת שנייה. אפשר לקבל אותה תוצאה גם אם גוזרים פעמיים את הפילטר G לקבלת הפילטר $d2G$ ואז מבצעים קונבולוציה על האות המקורי בעזרת $d2G$. הפילטר $d2G$ נקרא בשמות שונים וביניהם גם Laplacian.

הנה הדגמה:



החץ הירוק מסמן את "חציית האפס" (zero crossing), והוא המיקום של אמצע edge באות המקורי (אם מפצים על ההסטה שנגרמת במהלך ביצוע הקונבולוציה התלויה בגודל הפילטר. במקרה זה הפילטר באורך 400 ותיגרם הסטה באורך 200).

מסתבר שקל יותר למצוא מיקום מדויק של edges בעזרת חציית האפס (המושגת בעזרת החלקה ונגזרת שנייה) לעומת מציאת מקסימום (המושגת בעזרת החלקה ונגזרת ראשונה).

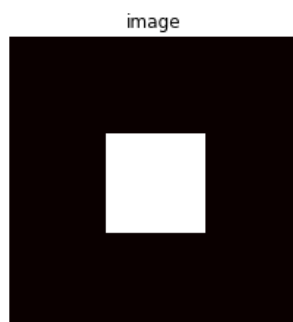
סיכום ביניים:

הדגמנו שימוש בקונבולוציה עם פילטרים שונים (נגזרת ראשונה, נגזרת שנייה, החלקה) ובצירופים שלהם לביצוע פעולות על אותות חד ממדים.

כעת נעבור לביצוע קונבולוציה דו ממדית (על תמונות \ מטריצות).

אפשר להשתמש בקונבולוציה של פילטר דו ממדי (מטריצה של ערכים) עם תמונה לביצוע פעולות מועילות.

נדגים מציאת edges בתמונה ללא רעש:



image

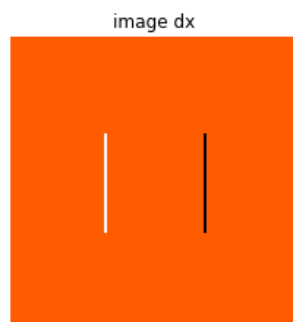


image dx

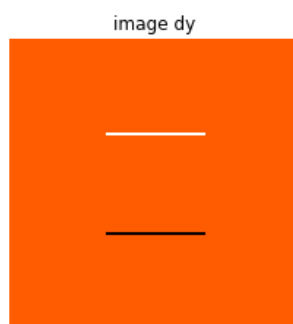
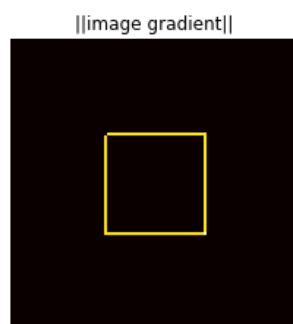


image dy



||image gradient||

לתמונה המקורית ריבוע לבן במרכז (לבן = 1, שחור = 0).

הפעלת נגזרת בציר ה x מגדישה edges אנכיים.

הפעלת נגזרת בציר ה y מגדישה edges אופקיים.

שימו לב לערכים של ה edges אחרי הגזירה.

לקבלת edges בעלי ערך חיובי בלבד אפשר לבצע פעולות שונות. כאן חושב גודל הגרדיאנט של התמונה.

גרדיאנט הוא וקטור המכיל את הנגזרות הכיווניות:

$$\nabla f = \left[\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right]$$

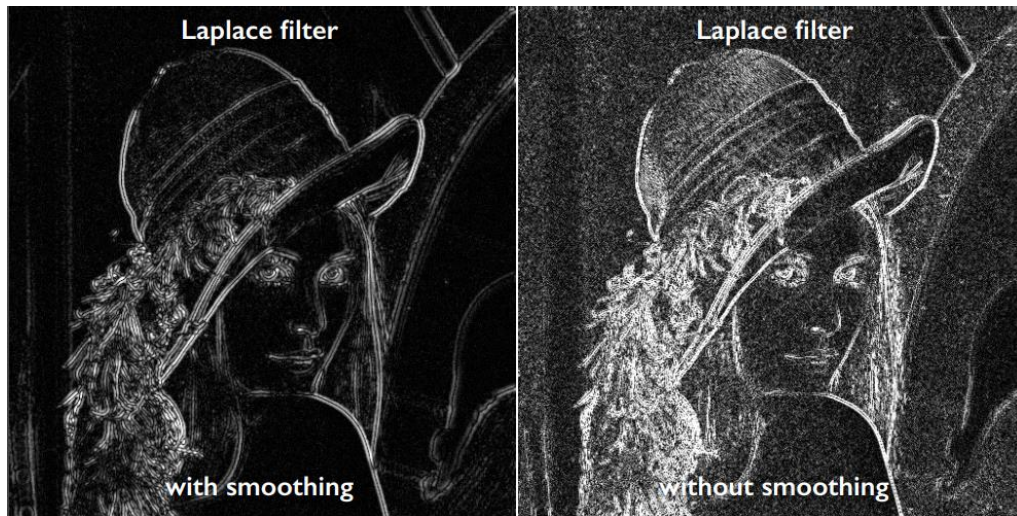
לחישוב גודלו נשתמש ב

$$\|\nabla f\| = \sqrt{\left(\frac{\partial f}{\partial x}\right)^2 + \left(\frac{\partial f}{\partial y}\right)^2}$$

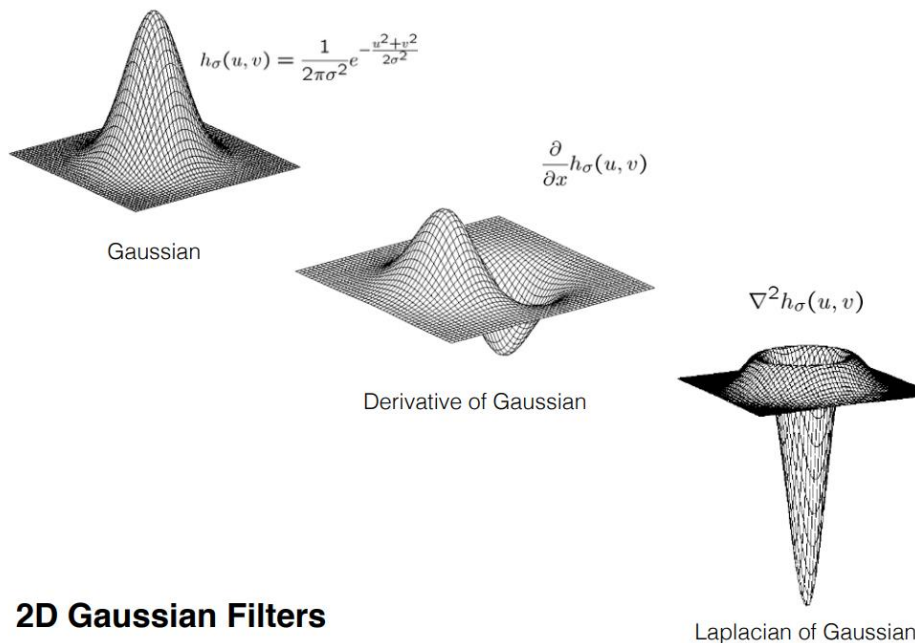
הנה פיסת קוד לביצוע חישוב זה:

```
gradient_norm = np.sqrt(np.square(image_dx) +  
np.square(image_dy))
```

גם בתמונות יש בד"כ רעש ועל כן משתמשים בשילובים שונים של החלקות וגזירות.
הנה התוצאה של הפעלת פילטר נגזרת שניה עם ובלי החלקה על תמונה



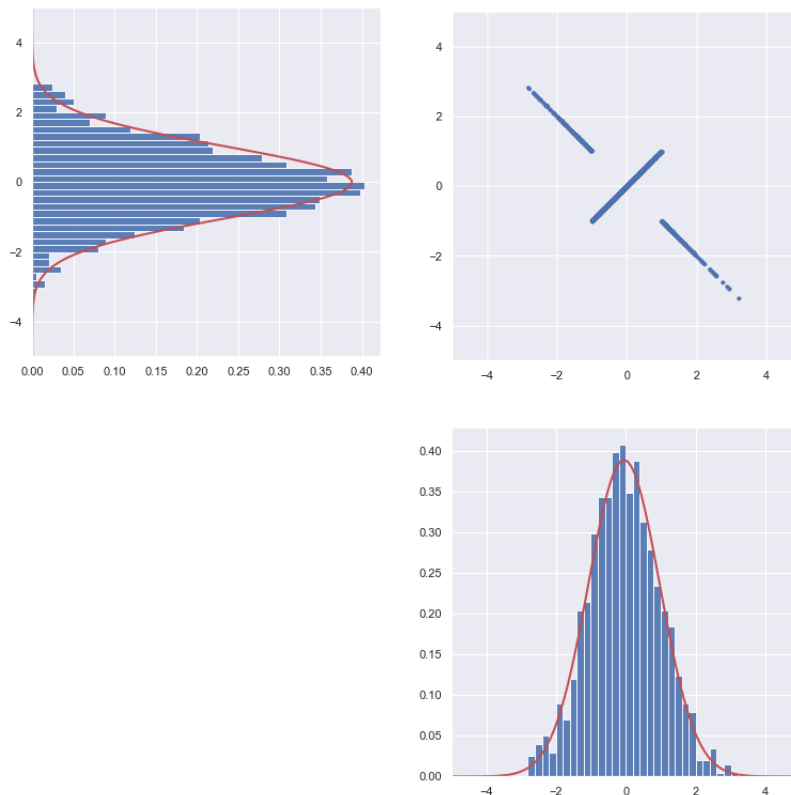
והנה הדגמה לצורת הפילטרים הדו ממדים כמשטחים (הערכים שבמטריצות מוצגים כאן כגובה המשטחים):



2D Gaussian Filters

שאלה 2

בשאלה זו תדגימו התפלגויות גאוסיאניות חד ממדיות והתפלגות משותפת שלהן. בשיעור הצגנו את המקרה הזה:



בו נראית התפלגות משותפת של שני משתנים (bivariate distribution) שאינה גאוסיאנית, בעוד שתי ההתפלגויות השוליות שלה (marginal distributions) הן גאוסיאניות.

א. כתבו קוד המשחזר את התוצאה המוצגת, על פי ההנחיות:

- לתוך x_1 הגרילו 1000 ערכים מהתפלגות נורמלית בעלת ממוצע 0 וסטיית תקן 1.
- x_2 יוגדר כפי שהוצג בשיעור.
- הציגו את הערכים x_1 כנגד x_2 בתצוגה דו ממדית. הקפידו על הצגה מרובעת ועל טווח ערכים מתאים לצירים.
- להצגת ההתפלגות השולית x_1 השתמשו בהיסטוגרמה של ערכי x עם 30 bins. הקפידו על טווח ערכים מתאים לציר האופקי.

על גבי ההיסטוגרמה נראית באדום עקומת ההתפלגות הגאוסיאנית המתאימה לערכי x_1 . אפשר היה להשתמש בהגדרת הממוצע וסטיית התקן המקוריים $(0,1)$, ולייצר את העקומה האדומה ולהציגה.

- במקום זאת השתמשו בהערכת הפרמטרים של התפלגות גאוסיאנית מתוך הנתונים x_1 . ההערכה תניב ממוצע וסטיית תקן ספציפיים למדגם.
- השתמשו בפרמטרים שהוערכו לדגימה של 100 ערכי y מ PDF גאוסיאנית כפונקציה של ווקטור x ובו 100 ערכים ברווחים קבועים בין -5 ל $+5$.
- הציגו את ה PDF כעקומה באדום על גבי ההיסטוגרמה. שימו לב שצריך לנרמל את ערכי העקומה להשגת ההתאמה להיסטוגרמה. לחלופין יש להשתמש בהיסטוגרמה מנורמלת. הסבירו בהערה בקוד את הנרמול בו אתם משתמשים.
- בדומה לטיפול ב x_1 בצעו את השלבים המתאימים עבור x_2 . שימו לב לסיבוב התצוגה ב 90 מעלות.

לפתרון סעיף זה עליכם למצוא בתיקוד של הספריות הרלוונטיות את הפונקציות המתאימות לביצוע:

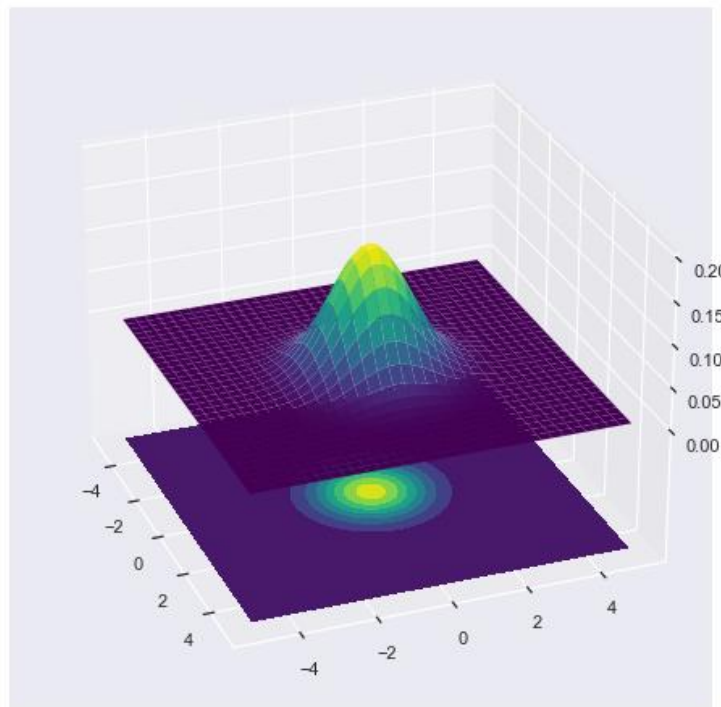
- הגרלה של משתנה מקרי המתפלג גאוסיאנית על פי ממוצע וסטיית תקן ידועים.
- הצגה של נקודות בדו ממד.
- שליטה בערכי הצירים של התצוגה, בגודלה וביחס הצירים כדי להציג תצוגות מרובעות.
- ייצור היסטוגרמה בעלת מספר bins רצוי.
- ביצוע התאמה (fit) של התפלגות גאוסיאנית לנתונים כך שישוערכו הממוצע וסטיית התקן של נתוני המדגם.
- ייצור של ערכי התפלגות גאוסיאנית בהינתן ממוצע וסטיית תקן ועבור טווח ערכי x ידוע.
- הצגה של עקומה רציפה על גבי גרף ההיסטוגרמה.
- הצגה של היסטוגרמה אנכית או אופקית ועקומה רציפה בהתאם.

ספריות ומחלקות מומלצות: `numpy`, `scipy.stats.norm`, `matplotlib (hist)`

1. להגשה: קובץ קוד בשם `ex2_a_1.py` המחשב ומציג את הגרפים.
2. העתיקו את הגרפים שהקוד שלכם מייצר למסמך הפתרון. נוח יותר לייצר שלושה גרפים נפרדים ולהעתיק אותם על פי הדוגמה למעלה לתוך מסמך הפתרון. אפשר גם לייצר גרף יחיד ובתוכו 3 subplots במקומות המתאימים לקבלת תצורה דומה.

ב. כעת תדגומו ייצור של התפלגות משותפת עבור שתי התפלגויות חד ממדיות בלתי תלויות, על פי ההנחיות:

- השתמשו בשתי ה PDF מן השאלה הקודמת (הערכים שהוצגו בעקומות האדומות). אלה הן דגימות של 100 ערכים מתוך התפלגויות גאוסיאנית.
- חשבו בעזרת ה PDF החד ממדיות התפלגות משותפת בהנחה שההתפלגויות החד ממדיות בלתי תלויות. התוצאה צריכה להיות מטריצה של 100 על 100 ערכים.
- הציגו את ההתפלגות המשותפת כך:



ספריות ומחלקות מומלצות:

matplotlib (plot_surface , contourf)

ב1. להגשה: קובץ קוד בשם ex2_b_1.py ובו החישובים וההצגה.

ב2. העתיקו את הגרף שהקוד שלכם מייצר למסמך הפתרון.

להרחבה והסבר נוסף על התפלגות נורמלית דו ממדית והצגותיה ראו כאן

[3D & Contour Plots of the Bivariate Normal Distribution – Data Science Genie](#)

שאלה 3

בשאלה זו תשתמשו בחוק בייז לניתוח הסתברותי.

נתונה מערכת החלטה רפואית. המערכת מזהה מחלה מסוימת על פי תוצאות של בדיקה. נזכיר כי תוצאה positive משמעה "חולה".

למערכת False Positive Rate (% זיהוי של בריא כחולה) של 5%, False Negative Rate (% זיהוי של חולה כבריא) של 5%. שכיחות המחלה באוכלוסייה הכללית היא 1 ל 100.

א. אדם נבדק בעזרת המערכת וקיבל תשובה חיובית (כלומר חולה). מה הסיכוי שהוא חולה? חשב והסבר.

ב. ידוע כי השגיאות בבדיקה אקראיות ובלתי תלויות. אותו אדם מסעיף א' ביצע את הבדיקה פעם נוספת וקיבל תשובה חיובית. מה הסיכוי שהוא חולה? הסבר.

שאלה 4

נתונות 6 דוגמאות קלט מעל הקו, 3 העליונות מניבות -1 ושלוש התחתונות מניבות +1: מה הניבוי לדוגמה שמתחת לקו?

A Learning puzzle

			$f = -1$
			$f = +1$
			$f = ?$

4א. רשמו 3 כללים שונים על פיהם הניבוי לדוגמה שמתחת לקו הוא -1.

4ב. רשמו 3 כללים שונים על פיהם הניבוי לדוגמה שמתחת לקו הוא +1.

שימו לב: כל הכללים צריכים לנבא נכונה את כל הדוגמאות שמעל הקו.

הגשה

- א. תאריך הגשה: עד יום ראשון, 27.11.22, בשעת חצות הלילה (המעבר ליום שני).
 - ב. ניתן להגיש בזוגות. אסור לעבוד בקבוצות גדולות יותר. **הגשה ב moodle**.
 - ג. יש לכתוב **שם \ שמות + ת"ז** בראשית כל מסמך מוגש (**כולל בקבצי הקוד**).
 - ד. כל מגיש (ביחיד או בזוג) צריך לדעת להסביר כל מה שנעשה בפתרון המוגש. חלק מן המגשים ידרשו להסביר את הפתרון שלהם למרצה.
 - ה. יש להגיש מסמך Word המכיל את כל התשובות לתרגיל. שם מסמך זה יהיה **ex1.docx**.
- הקפידו שמספור סעיפי התשובות שלכם יהיה **זהה** למספור סעיפי השאלות.
- ו. לכל פונקציה בקוד צריך להיות תיעוד במתכונת הזו

```
def add(a, b):  
    """  
    Sum up two integers  
    Arguments:  
    a: an integer  
    b: an integer  
    Returns:  
    The sum of the two integer arguments  
    """  
    return a + b
```

- ז. יש להגיש את כל הקוד לתרגיל בקבצים על פי השמות שניתנו למעלה.
- ח. כל הקבצים ישכנו בתוך **קובץ דחוס** (zip) הכולל את שמכם.
שם הקובץ למגיש יחיד:
EX1Family1Name1
שם הקובץ לשני מגישים:
EX1Family1Family2