

## אלגוריתם Gradient Decent וחוק עדכון

### המשקולות:

השגיאה היא פונקציה תלויה משקולות. על מנת למזער אותה עלינו להגיע לנקודת המינימום שלה וכך נקבל את המשקולות שמביאים את המינימום. 1. התחל ממשקולות אקראיים קטנים ובצע epochs שוב ושוב עד לקיום תנאי עצירה (השגיאה לא מפסיקה לרדת או מסי' איטרציות):  
1.1 בכל epoch עדכן את המשקולות באופן הבא:

$$\Delta w_i = -\lambda \frac{\partial \text{loss}(w)}{\partial w_i}$$
$$w_i = w_i + \Delta w_i$$

היפר פרמטר, קבוע הלמידה  $\lambda$ . קיים גם בגרסת online שבו עדכון המשקולות אחרי כל דוגמא.

### גרסה לינארית:

מודל פרמטרי לשיערוך היפותזה בעזרת משערוך לינארי (קו או מישור ישר).

עבור  $n$  מאפיינים, ההיפותזה נראת כך:

$$y = wx = w_0 \cdot 1 + w_1 \cdot x_1 + \dots + w_n x_n$$

השגיאה נראת כך:

$$\text{loss} = \text{MSE} = \frac{1}{n} \sum_{i \in D} (t_i - y_i)^2$$

### מציאת המשקולות:

ישנן 2 דרכים למצוא את המשקולות כך שימזערו את השגיאה:

1. המשוואה הנורמלית  $w = (x^T x)^{-1} x^T t$

שיטה מדויקת אך מתאימה למימדים קטנים בלבד

2. אלגוריתם GD כפי שתואר מעלה. לשם כך, הגזירה היא

$$\frac{\partial \text{MSE}}{\partial w} = \frac{\partial \text{MSE}}{\partial y} \frac{\partial y}{\partial w} = \frac{1}{n} \sum_{i \in D} (t_i - y_i) \cdot x_i$$

אם סט האימונים קטן נעדיף את המשוואה הנורמלית, אחרת את אלגוריתם GD.

### נרמול שדות:

פונקצית השגיאה היא קערה, על מנת לוודא שהיא סימטרית יש לנרמל את ערכי המאפיינים:

1.  $\frac{x - \min}{\text{range}}$  = Min max scaling

2.  $x - \text{mean}(x)$  = Mean normalization

3.  $\frac{x}{\text{sd}(x)}$  = SD scaling

### גרסה פולינומאלית:

הוספת מאפיינים שהם טרנספורמציות של מאפיינים קיימים (בודד או מורכב). התאמת ההיפותזה = נוסף מאפיינים לפי דרגת הפולינום שבו מתעניינים להיפותזה (ניתן להתאים פולינומים מכל דרגה) על ידי הוספת כלל האפשרויות המתאימות. לכל מאפיין חדש משקל חדש.

## גרסה לוגיסטית: (קלסיפיקציה)

מודל פרמטרי לסיווג קטגוריות (קלסיפיקציה) בעזרת קלסיפייר (יוצר גבול הפרדה בין קטגוריות). עבור  $n$  מאפיינים, ההיפותזה נראת כך:

$$z = wx = w_0 \cdot 1 + w_1 \cdot x_1 + \dots + w_n x_n$$
$$y = g(z) = \frac{1}{1 + e^{-z}}$$

השגיאה נראת כך:

$$\text{loss} = \text{CE} = -t \cdot \log(y) - (1 - t) \log(1 - y)$$

$$C(y, t) = \begin{cases} -\log(y) & \text{if } t = 1 \\ -\log(1 - y) & \text{if } t = 0 \end{cases}$$

### מציאת המשקולות:

1. אלגוריתם GD כפי שתואר מעלה. לשם כך, הגזירה היא

$$\frac{\partial \text{CE}}{\partial w} = \frac{\partial C}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial w} = \frac{1}{n} \sum_{i \in D} (t_i - y_i) \cdot x_i$$

סיווג שאינו בינארי בעזרת קלסיפייר בינארי:

גישה one vs rest נבנה קלסיפייר לכל קטגוריה, כלומר נעביר קו חוצה בינה לבין הקטגוריות

האחרות, סה"כ  $k$  מסווגים. אימון: נלמד להפריד קטגוריה אחת מכל השאר. חיזוי: בהינתן דוגמה לסיווג נבדוק תחזית של כל מסווג ונחזיר את הקטגוריה שלה חוצה הסתברות הגבוהה ביותר.

גישה one vs one נבנה קלסיפייר לכל זוג קטגוריות, כלומר נעביר קו חוצה בין כל זוג קטגוריות. סה"כ

$C(k, 2)$  מסווגים. אימון: נלמד להפריד קטגוריה אחת מכל השאר. חיזוי: מספר אפשרויות. אחת תהיה הצבעת הרוב, אחרת תהיה חישוב הסתברות ממוצעת לכל קטגוריה ובחירת קטגוריה שמקבלת את ההסתברות המקסימלית.

$$\text{MCCE}_D(y, t) = -\frac{1}{m} \sum_p \sum_i t_{pi} \log(y_{pi}) + (1 - t_{pi}) \log(1 - y_{pi})$$

### גמישות מודלים:

גמישות של מודל מתבטאת בדרגת הפולינום שלו. גמיש יותר ככל שהדרגה עולה. **Under fitting** = אין התאמה מספיק טובה בין המודל הנוצר לנתונים. **Over fitting** = התאמה מידי טובה בין המודל הנוצר לנתונים. על מנת להגמיש את המודל ניתן: להוסיף מאפיינים, להשתמש בהיפותזות מורכבות וגמישות יותר (רשתות). **הגמשה עבור גרסיה:** שיטת **LWLR**. נקודות יותר קרובות ישפיעו יותר מנקודות רחוקות. קרבה מחושבת בעזרת

$$\beta_i = e^{-\frac{\|x_i - x\|^2}{2\tau^2}}$$

ואז פונקצית השגיאה היא

$$\text{WMSE}_w = \frac{1}{2m} \sum_i \beta_i (wx_i - t_i)^2$$

$\tau$  קטן יותר המשקל של נקודות שהן מרוחקות קטן ואילו של נקודה קרובה גדל, ככל שיותר קטן יותר  $\tau$ . over fitting גדול משמע גרסיה לינארית רגילה

<p>= Lasso regularization (L1)</p> <p>נעניש את הפונקציית השגיאה על ידי הוספת ריבוע הנורמה מסדר ראשון <math>\ w\ _1 = \sqrt{\sum_{i=1}^n  w_i }</math> ואז, נוסחת השגיאה המעונשת תהיה</p> <p><math>regMSE(w) = MSE(w) + \gamma \sum_{j=1}^n  w_j </math></p> <p>אם כך, כלל עדכון המשקולות משתנה בהתאם:</p> $\Delta w = -\lambda \frac{\partial loss}{\partial w_i} - \gamma$ <p>שיטה זו מסוגלת לבצע בחירת מאפיינים (לאפס).</p> <p>= Elastic net regularization</p> <p>שילוב בין שתי הנורמות. במשקולות הקטנים לאסו משפיע ואילו במשקולות הגבוהים, הרידג' (אבל הלאסו ממתן קצת). הנוסחה</p> $\frac{lasso+ridge}{2}$ <p><b>3. אנסמבלים:</b></p> <p>אנסמבלים הם צירופים של הרבה מודלים שונים, כאשר מתוכם יבחר המודל הטוב ביותר. יצירה:</p> <p>= Bagging</p> <p>יצירת אנסמבלים על ידי דגימות שונות מתוך סט האימונים. מאמנים מודלים שונים על תת קבוצות שונות של הנתונים, ואת הפלט ממזעים או החלטת הרוב.</p> <p>(1) שיטת <b>boot strap</b> = שיטת דגימה בה מייצרים הרבה קבוצות אימון שונות ע"י דגימה אקראית עם חזרות מתוך קבוצת אימון יחידה. לכל קבוצה מותאמת קבוצת וולידציה הכוללת דוגמאות שלא נלקחו לדגימה.</p> <p>(2) שיטת <b>k fold</b> = מיצרים k קבוצות אימון ע"י שמוציאים בכל פעם 1/k מהנתונים כוולידציה, ומשאירים בקבוצת האימון את היתר. קבוצות הוולידציה הן זרות זו לזו.</p> <p>שיטת <b>boosting</b> = אנסמבלים של מודלים מאומנים על קבוצות אימון שונות כך שכל מודל מתמקד בדוגמאות שהקודמים לו שגו בהן. משתמשים בסדרה של מודלים "חלשים" (עם שגיאת ביאס גבוהה) וכל מודל מתרכז בלמידה של המקרים שהמודלים האחרים לא הצליחו לסווג נכון. משתמשים באותם נתונים, אך מגדילים את משקל הדוגמה (תדירות הופעתה) במידה והמודלים הקודמים לא סיווגו אותה נכון. מקבלים סידרת מומחים (שכל אחד "חלש" מידי – ולא יכול לעשות Overfitting) שנותנים תחזיות שונות למיקרים קשים. את המודלים השונים שמתקבלים ניתן לקבל ע"י למשל הצבעת הרוב או מיצוע (גיאומטרי) של הסתברויות – (אפשר לשקלל מודל על פי חשיבות השגיאות שמטופלות על ידו) או על פי מידת הצלחתו על וולידציה.</p> <p><b>פירוק שגיאת MSE:</b></p> $loss = (Bias)^2 + Variance + irreducible$ <p>בד"כ שגיאת ביאס משמע under fitting ואילו שגיאת שונות משמע over fitting. ככל שמעלים גמישות של מודל הביאס יורד והשונות עולה.</p>	<p><b>הקשחת המודלים:</b></p> <p><b>1. Feature selection</b></p> <p>לפי עקרון התער על אוקאם, " אין להעמיד ריבוי ללא צורך", כלומר העדפת ההסבר הפשוט יותר על פני המורכב. במקרה הזה, העדפה למודלים בעלי פחות פארמטרים. אלגוריתם בחירת פרמטרים יחזיר רק את המאפיינים הטובים ביותר.</p> <p>=Forward feature selection</p> <div style="border: 1px solid black; padding: 5px;"> <p>1) Start with empty feature set</p> <p>2) For <math>k = 1, \dots, n</math>:</p> <p>2.1) For all features <math>f_k</math> not used by model:</p> <p>2.1.1) Try adding each of the missing features, train and save the loss</p> <p>2.2.2) Create model <math>M_k</math> by adding the feature that provides the best loss</p> <p>3) Choose the best <math>M_k</math> using cross validation</p> </div> <p>יווצרו <math>n^2</math> מודלים (סדרה חשבונית), <math>n</math> שיערוכי שגיאה cross validation.</p> <p>=Backwards feature selection</p> <div style="border: 1px solid black; padding: 5px;"> <p>1) Start with full feature set</p> <p>2) For <math>k = 1, \dots, n</math>:</p> <p>2.1) For all features <math>f_k</math> not used by model:</p> <p>2.1.1) Try removing each of the missing features, train and save the loss</p> <p>2.2.2) Create model <math>M_k</math> by removing the feature that provides the best loss</p> <p>3) Choose the best <math>M_k</math> using cross validation</p> </div> <p>יווצרו <math>n^2</math> מודלים (סדרה חשבונית), <math>n</math> שיערוכי שגיאה cross validation.</p> <p>=Hybrid feature selection</p> <div style="border: 1px solid black; padding: 5px;"> <p>1) Start with full feature set</p> <p>2) For <math>k = 1, \dots, n</math>:</p> <p>2.1) add a feature in Forward selection</p> <p>2.2) remove a feature in Backward selection</p> <p>2.3) save <math>M_k</math></p> <p>3) Choose the best <math>M_k</math> using cross validation</p> </div> <p>יווצרו <math>n^2</math> מודלים (סדרה חשבונית), <math>n</math> שיערוכי שגיאה cross validation.</p> <p><b>2. רגולריזציה:</b></p> <p>הענשת משקולות. נגדיר <math>\gamma</math> קבוע רגולריזציה, ונוסחת שגיאה חדשה, מעונשת.</p> <p>= Ridge regularization (L2)</p> <p>נעניש את הפונקציית השגיאה על ידי הוספת הריבוע של הנורמה מסדר שני (שהם ריבועים): <math>\ w\ _2 = \sqrt{\sum_{i=1}^n w_i^2}</math> ואז, נוסחת השגיאה המעונשת תהיה</p> $regMSE(w) = MSE(w) + \frac{\gamma}{2} \sum_{j=1}^n w_j^2$ <p>אם כך, כלל עדכון המשקולות משתנה בהתאם:</p> $\Delta w = -\lambda \frac{\partial loss}{\partial w_i} - \gamma w_i$
--	--

### 1. K-fold CV

חלק את הנתונים המותוגים ל  $K$  קבוצות. שמור  $\frac{1}{K}$  מהנתונים כקבוצת ולידציה כל פעם. אמן על  $\frac{K-1}{K}$  מהנתונים ובדוק על הוולידציה.

בצע  $K$  אימונים שונים בכל פעם על קבוצת ולידציה אחרת. סה"כ  $K$  מודלים יחושבו.

### 2. Leave K out CV

חלק את הנתונים המותוגים ל  $K$  קבוצות. קבוצה אחת כקבוצת ולידציה כל פעם. אמן על  $D - K$  מהנתונים ובדוק על הוולידציה. בצע

$C\left(\frac{D}{K}\right)$  אימונים שונים בכל פעם על קבוצת ולידציה אחרת. סה"כ  $C\left(\frac{D}{K}\right)$  מודלים יחושבו.

### 3. Bootstrap Method

יכולות החיזוי והערכת השגיאה היו משתפרות אם היינו יכולים לייצר כמויות גדולות של נתונים מותוגים שאיתם היינו מייצרים הרבה קבוצות אימון. אבל, ברוב המקרים אין נתונים רבים. בשיטה הזו, מייצרים הרבה קבוצות אימון מקבוצת אימון אחת שאותה דוגמים רנדומית עם חזרות. קיים סיכוי

$\frac{1}{3}$  שחלק מהדוגמאות בסט לא יופיעו בקבוצת האימון, אותן ניקח כולידציה

### אלגוריתמים שונים:

#### 1. K nearest neighbors

**עבור קלסיפיקציה:** בהינתן אוסף דוגמאות ודוגמת מבחן חדשה, נשערך את הסתברות השייכות של הדוגמא לכל אחת מהקטגוריות על פי החלטת הרוב של  $k$  השכנים הקרובים ביותר.

**עבור רגרסיה:** בהינתן אוסף דוגמאות ודוגמת מבחן חדשה, נשערך את החיזוי של הדוגמא לפי ממוצע הערכים של  $k$  השכנים הקרובים ביותר. הגמישות יורדת ככל שיותר שכנים משתתפים בהחלטה. יש להחליט מי קרוב:

**עבור מאפיינים נומרים:**

✓ מרחק מנהטן l1 norm

$$d(A, B) = |x_A - x_B| + |y_A - y_B|$$

✓ מרחק אוקלידי l2 norm

$$d(A, B) = \sqrt{(x_A - x_B)^2 + (y_A - y_B)^2}$$

✓ נורמת האינסוף

$$d(A, B) = \max |A_i| - \max |B_i|$$

**עבור מאפיינים שאינם נומרים:**

✓ מרחק קוסינוס שמודד קרבה לפי דמיון

$$d(A, B) = \frac{A \cdot B}{|A||B|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

שני וקטורים מאונכים אחד לשני הדמיון שלהם יהיה אפס, ואילו וקטורים עם אותה זווית, דמיון אחד.

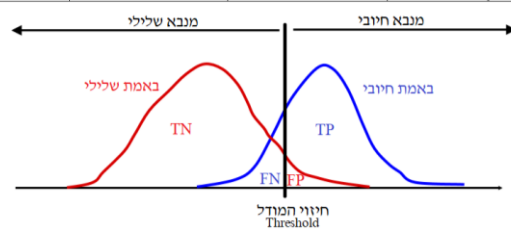
### מדדת הצלחת המודלים:

נשתמש  $KPI's$  = Key Performance Indicators

Recall, TPR	$\frac{TP}{TP + FN} = 1 - FNR$
Specificity, TNR	$\frac{TN}{TN + FP} = 1 - FPR$
Precision / PPV	$\frac{TP}{TP + FP} = 1 - FDR$
false negative rate	$\frac{FN}{FN + TP} = 1 - TPR$
false positive rate	$\frac{FP}{FP + TN} = 1 - TNR$
false discovery rate	$\frac{FP}{FP + TP} = 1 - PPV$
Accuracy מדד כללי טוב להצלחת מודל	$\frac{TP + TN}{TP + TN + FP + FN}$
F1 score עוזר כשחלק מהקטגוריות נדירות	$2 \cdot \frac{P \cdot R}{P + R}$
Balanced accuracy	עבור מודל שאינו בינארי הוא הממוצע של recall, אחרת, ממוצע בין recall לבין הספציפיות

### ניתן לעזור במטריצת עמימות

Confusion matrix		החיזויים (y)		
		חיזוי חיובי	חיזוי שלילי	
הסיווג האמיתי (t)	באמת חיובי	True Positive	False Negative שגיאה מסוג 2	TPR or recall $\frac{TP}{P}$
	באמת שלילי	False Positive שגיאה מסוג 1	True Negative	TNR = $\frac{TN}{N}$
		PPV or Precision $\frac{TP}{\text{Predicted P}}$	NPV = $\frac{TN}{\text{Predicted N}}$	Accuracy $\frac{TP + TN}{\text{ALL Examples}}$



נוכל ליצור מטריצת עמימות גם עבור קלסיפיקציה לא בינארית. גודל המטריצה יהיה כמספר הקטגוריות. ישנן  $n$  קטגוריות אז גודל המטריצה יהיה  $n \times n$ . ואז יהיה ניתן לחשב את כל המדדים שנלמדו על כל קטגוריה בנפרד.

### שיעור הצלחה בעזרת ולידציה CV 33

שגיאת הולידציה עבור  $k$  דוגמאות היא

$$CV_{loss} = \frac{1}{k} \sum_v loss_v$$

$$p(v_1, v_2, \dots, v_n | k) \approx \prod_{j=1}^n p(x_j = v_j | k)$$

✓ סדר ההופעה אינו חשוב

$$p(a_i = w_j | k) = p(a_m = w_j | k)$$

אלגוריתם בייס הנאיבי:

בהתאם להנחות הנאיביות, נבנה האלגוריתם

Naïve Bayes Classifier (D, K, x):

For each prediction category  $k=1 \dots K$ :

- $P(k) \approx \frac{n_k}{n}$
- For each value  $v_j$  of each attribute  $x_j$ :

$$P(v_j | k) = p(x_j = v_j | k) \approx \frac{n_{v_j \cap k}}{n_k}$$

Return

$$\text{predict} = \operatorname{argmax}_k \{P(k) \cdot \prod_{j=1}^n P(v_j | k)\}$$

ישנו צורך לבצע תיקונים לאלגוריתם עבור

קטגוריות נדירות **m-estimate**:

כאשר אין מספיק דוגמאות מקטגוריה  $k$  ושבהן גם

$x_i = v_i$ , מוסיפים  $m$  דוגמאות וירטואליות

מקטגוריה  $k$  שמתוכם  $P(x_i = v_i)$  הם בעלי ערך

$x_i = v_i$ .  $P(x_i = v_i)$  הינו השערוך prior של  $v_i$ .

$m$  הינו היפר פאראמטר שמצפה על גודל המדגם או על התלויים.

$$\hat{P}(x_i = v_i | k) \approx \frac{n_{v_i \cap k} + m \cdot P(x_i = v_i)}{n_k + m}$$

אם לא ניתן לשערך את ההסתברות הקודמת

$P(x_i = v_i)$  נשתמש בתיקון **Laplace smoothing**:

אפשר להניח התפלגות אחידה בין הערכים  $v_1 \dots v_n$

ולשערך  $P(x_i = v_i) \approx \frac{1}{n}$ . ואז ההסתברות תהיה

$$\hat{P}(x_i = v_i | k) \approx \frac{n_{v_i \cap k} + m \cdot \frac{1}{|v|}}{n_k + m}$$

בפרט, אם נוסיף מספר דוגמאות כמספר הערכים

של הקטגוריה החסרה  $m = |v|$  נקבל  $\frac{n_{v_i \cap k} + 1}{n_k + |v|}$ .

### 3. עצי החלטה ויער רנדומי

גישה שבה ההיפותזה היא עץ החלטה. נתונה קבוצת

אימון עם דוגמאות ממספר קטגוריות. נחפש שאלה

על אחד המאפיינים שתפצל את קבוצת האימון

לשתיים (או יותר) קבוצות הומוגניות ככל האפשר.

ניצור צומת הכולל את השאלה שנבחרה ונפצל את

קבוצת הדוגמאות לקבוצות (ילדים בעץ) בהתאם

לשאלה. בצורה רקורסיבית, נמשיך לפצל את כל

קבוצות הצאצאים עד שנצליח ליצור עלים

הומוגניים (הפרדה מלאה) או שנעצור קודם על פי

קריטריון אחר.

✓ מרחק קורלצית פירסון שמוזדד קרבה לפי קורלציה

$$d(A, B) = \frac{\sum_{i=1}^n (A_i - \bar{A})(B_i - \bar{B})}{\sqrt{\sum_{i=1}^n (A_i - \bar{A})^2} \sqrt{\sum_{i=1}^n (B_i - \bar{B})^2}}$$

$k$  הינו היפר פארמטר עבור האלגוריתם.

שכלול האלגוריתם = **weighted k-NN**:

הבאת משקל רב יותר לשכנים הקרובים ביותר.

בעזרת דימיון גאומטרי, מאפשרת שליטה בגמישות:

$$w_i = e^{-\frac{\|x_i - x\|^2}{2\tau^2}}$$

חסרון האלגוריתם:

רגיש מאוד למספר המימדים. הדמיון יכול להטעות.

ניתן לפתור על ידי הורדת מימדים.

## 2. למידה בייזיאנית

התבססות על תורת ההסתברות, בפרט על משפט

בייס וכן על 2 הנחות נאיביות.

$$P(B|A) = \frac{P(A \cap B)}{P(A)} \quad \text{נוסחת ההסתברות המותנית:}$$

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} \quad \text{חוק בייס:}$$

נוסחת ההסתברות השלמה:

$$P(B) = \sum_i^k P(B|A_i)P(A_i) \quad \text{ומתוכה ניתן לשנות}$$

$$P(A_i|B) = \frac{P(B|A_i) \cdot P(A_i)}{\sum_i^k P(B|A_i)P(A_i)} \quad \text{את נוסחת בייס:}$$

הסתברות מותנית נקראת **posterior** ואילו

הסתברות של מאורע כלשהו נקראת **prior**.

עבור המאורע  $x|k$  נרצה לשערך את הייתכנות שלו.

ישנן 2 דרכים.

הערכה מקסימלית להסתברות **MAP = posterior**:

$$y_{MAP} = \operatorname{argmax}_k \{P(x|k) \cdot P(k)\}$$

כלומר, נחשב עבור כל ערכי  $x$  את ההערכה ונבחר

את המקסימלית. היא השערוך להסתברות הרצויה.

הערכה מקסימלית לסבירות **ML = likelihood**:

$$y_{ML} = \operatorname{argmax}_k \{P(x|k)\}$$

כלומר, נחשב עבור כל ערכי  $x$  את ההערכה ונבחר

את המקסימלית. היא השערוך להסתברות הרצויה.

נשים לב כי אילו הערכות בלבד שאינן מחשבות את

ההסתברות עצמה. על מנת לחשב את ההסתברות

עצמה יש לחשב את כל הנדרש כדי לקבל:

$$p(k|x) = \frac{k \cap x}{p(x)}$$

הנחות בייס:

✓ המאפיינים אינם תלויים זה בזה בהינתן

קטגוריה  $k$  ולכן ההסתברות לוקטור קלט

כלשהו הוא כפל ההסתברויות של כל

המאפיינים בו

נסדר את ערכי הערכי המאפיין בסדר עולה (וברשימה נפרדת את ערכי התוויות המתאימים). נבדוק את רווח המידע  $\text{info gain}$  (=ממוצע הערכים במעבר) רק בנקודות שבהם מתחלפת התווית. מובטח לרווח המידע המקסימלי להיות באחד המעברים

Temp:	40	48	60	72	80	90
PlayTennis:	No	No	Ye	Ye	Ye	No

ניתן גם לאפשר רגרסיה בעץ: ממוצע התוויות של הדוגמאות הנופלות בקבוצה (צומת) נלקח כחיזוי  $\hat{y}$  של הקבוצה. כל עלה, לכן, חוזר את הערך הנומרי הממוצע של הדוגמאות שבעלה.

$\text{sum squared error}$  של הדוגמאות בקבוצה נלקח כמדד לשיפור (במקום אנטרופיה). כלומר, סכום ריבועי ההפרשים בין התוויות לממוצע של הדוגמאות בצומת יהיה מדד השיפור. הרווח של כל פיצול מחושב ע"פ השיפור במדד  $SSE$  בעקבות הפיצול.

**מדד Genie** = ממד זה דומה למדד האנטרופיה וגם הוא מודד אי סדר. עבור משתנה אקראי בינארי, המדד יחושב כך

$$G = \sum_{i=1}^N p_i \cdot (1 - p_i)$$

הומוגיות,  $G = H = 0$ . בשונה ממד האנטרופיה, כאשר קבוצה היא הטרוגנית,  $\max G = \frac{1}{2}$ . כאשר הסתברויות לקטגוריות קרובות לאחד או לאפס, שני המדדים יתנו מספר נמוך (קרוב ל 0 מלמעלה). שני המדדים מקבלים ערך גדול מאפס כאשר יש הרבה קטגוריות בהסתברות זהה, אולם  $G$  חסום ב 1 ואילו  $H$  יכול להגיע למספרים גבוהים. עצים יכולים לבצע  $\text{overfitting}$  ולהתייחס למאפיינים לא חשובים ולכן עדיף להשתמש באנסמבל של עצים.

Bagging using bootstrap:

נבצע דגימה עם חזרות כדי לגדל  $B$  עצים שונים המורכבים מ  $B$  מדגמים שונים קצת זה מזה. נעזר בתוצאות החיזוי של העצים השונים: ברגרסיה ממצעים את החיזוי, בקלסיפיקציה משתמשים בהחלטת הרוב או מיצוץ ההסתברות של העלה אליו הגענו בכל אחד מהעצים. חסרון בשיטה זו הוא שהעצים דומים מידי אחד לשני ולכן שגיאת השונות אינה יורדת משמעותית.

Random forest:

נרצה לגדל עצים הנבנים מקבוצות שונות של מאפיינים. נבנה עץ לכל מדגם בשיטת bootstrap. בכל פעם שבחרים צומת לפיצול, נגביל את המאפיינים לבדיקה ל  $m' < m$  מאפיינים הנבחרים רנדומית מתוך  $m$  המאפיינים שלרשותנו. כלומר, אם  $m' \ll m$ ,

מאפיין הינו קטגורי, ניתן לפצל בשתי דרכים:

- ✓ פיצול שהוא יהיה כל הקטגוריות האפשריות
  - ✓ פיצול שהוא בינארי על ידי בחירת תת קבוצה של ערכים, למשל, קטגוריה אחת מול כל השאר או לקבוצת ערכים כנגד הקבוצה המשלימה
- מאפיין נומרי, ניתן לפצל, גם כן, בשתי דרכים:

- ✓ ניתן לפצל אותו באופן בינארי על ידי בחירת סף  $x_i < \text{threshold}$
- ✓ ניתן להגדיר טווחי ערכים (Bins) ולפצל את העץ כאילו היו קטגוריות בדידות

מרחב ההיפותזות של כל העצים האפשריים הוא עצום. בהנחה שישנם  $n$  משתנים בוליאנים ושהעצים הם

מלאים, מספר העצים הוא  $2^{2^n}$ . על כן, נחפש באופן חמדני את העץ הטוב ביותר.

אם הקבוצה הומוגנית, ניצור עלה וסיימנו. אחרת, נבחר מאפיין שיפצל הכי טוב טוב לשתי קבוצות נפרדות. באופן רקורסיבי נמשיך לפצל עוד ועוד בנים עד שהקבוצות יהיו הומוגניות או שלא ניתן לפצל יותר. כדי לחשב מה הוא מאפיין הפיצול הטוב ביותר, נעזר **במדד האנטרופיה** = מדד להומוגניות, לסדר של קבוצה ביחס לקטגוריות המופיעות בה. ככל שיש יותר ערכים בקבוצה, אי הסדר גדל,

$$H(V) = \sum_{v=0}^1 -P(V=v) \log P(V=v)$$

כדי להמשיך ולאפיין, נשכלל את האנטרופיה

$$H_{tot} = \sum_{i=1}^{\#children} P(child_i) \cdot H(child_i)$$

כלומר, האנטרופיה **המשוכללת** היא סכום האנטרופיות של הבנים במכפלת הסיכוי לפיצול הזה על העץ. ועכשיו, נוכל לדעת כמה מידע הרווחנו בעזרת פיצול זה על ידי

$$IG = H(\text{parent}) - H_{tot} \quad \text{information gain}$$

להלן האלגוריתם החמדני:

Grow Tree( $S$ ):

if ( $t=0$ ) for all examples, return (new leaf(0))

if ( $t=1$ ) for all examples, return (new leaf(1))

Else

choose "best" feature  $x_j$  not used yet by

checking entropy and IG

$$S_0 = \text{all } \langle x, t \rangle \in S \text{ with } x_j = 0$$

$$S_1 = \text{all } \langle x, t \rangle \in S \text{ with } x_j = 1$$

return (new node ( $x_j < 0.5$ , Grow Tree( $S_0$ ),

Grow Tree( $S_1$ )))

נשים לב כי קיימת בעייתיות עבור מאפיינים שהם נומרים רציפים, למשל טמפרטורה. במקרה כזה, נחפש סף שלפיו ניתן לפצל. נבדוק את רווח המידע עבור מספר מקומות סף שבהם ניתן לפצל פיצול בינארי.



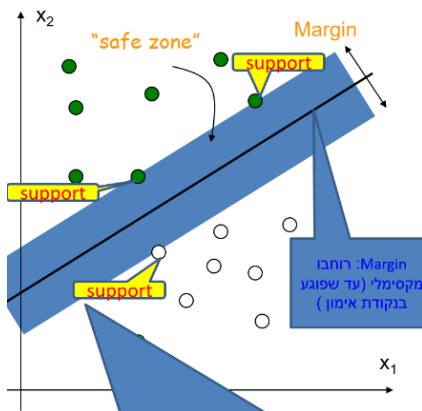
דוגמאות תוכלנה להופיע בצד הלא נכון של השוליים מבלי להצר אותו וכן, מיעוט דוגמאות תוכלנה אפילו להופיע בצד הלא נכון של מישור ההפרדה מבלי שישפיעו על מישור ההפרדה ואפילו על השוליים, משתמע מכך שההפרדה אינה לינארית מושלמת. רוצים למקסם את השוליים אבל עבור שוליים רחבים מוכנים מוכנים "לשלם" שמיעוט נקודות יכנסו לתוך הצינור ואפילו יסווגו כשגויים. הפתרון הוא להוסיף משתנים לבעית האופטימיזציה. נוסיף משתנה רפיון  $Slack\ variable = \epsilon_i$  לכל דוגמה,

אשר מאפשר לה להופיע בצד הלא נכון, כאשר  $\epsilon_i$  הוא התשלום שמשלים עבור חריגה מהשוליים. כעת, נרצה למעור את סכום התשלומים.

כאשר  $\epsilon_i = 0$  הדוגמה, תהיה מחוץ לשוליים ותסווג נכון. כאשר  $0 < \epsilon_i \leq 1$  הדוגמה תהיה בצד הלא נכון של השוליים אבל בצד הנכון של המישור ותסווג נכון. כאשר  $\epsilon_i > 1$  הדוגמה תהיה בצד הלא נכון של המישור המפריד ותסווג שגוי.

יהיה לנו תקציב  $C$  שאנחנו מוכנים לשלם עבור הנקודות החורגות הללו, ככל שהחריגה גדולה יותר, התשלום  $\epsilon_i$  גדול יותר. נרצה להרחיב את השוליים כמה שיותר, אך ובשום פנים לא לחרוג מהתקציב.  $C$  הינו היפר פרמטר שולט Bias – Variance tradeoff

לא יתאפשרו יותר מאשר  $C$  חריגות מהמישור המפריד. תקציב גדול, מאפשר להרבה נקודות לחרוג מהשוליים וכל נקודה שחורגת היא וקטור תמך. כשיש הרבה וקטורי תמך יש פחות רגישות לתזוזות בנקודות אלו, או להוספת נקודות תמך חדשות (שגיאית שונות קטנה יותר). מצד שני, תקציב גדול מאפשר סיווג לא נכון של דוגמאות אימון ולכן עלול להעלות את שגיאת ביאס.



בכל פעם שבחרים צומת לפיצול, האלגוריתם לא מורשה להסתכל על רוב המאפיינים. העצים שנגדל יהיו שונים, גם בגלל שנוצרו מקבוצות אימון שונות במקצת ובעיקר בגלל שאולצו להשתמש במאפיינים שונים.

בהינתן קבוצת דוגמאות  $S$ , רשימת מאפיינים מותרים והיפר פרמטר  $m'$ , להלן האלגוריתם מגבלה אקראית על מאפיינים הניתנים לבחירה:

Choose Best Feature ( $S$ , Feature set,  $m'$ ):  
 1) If Feature set is empty return (Null, Null)  
 2) Randomly select  $m'$  features from Feature set (if  $|Feature\ set| < m'$ , select all)  
 3) Find the Best Feature with the best IG  
 3.1) remove from Feature set, return (Best Feature, Feature set)

להלן האלגוריתם המלא לבניית עץ בינארי אקראי:

Grow Rand Best Tree ( $S$ , Feature set,  $m'$ ):  
 If  $t=0$  for all  $\langle x, t \rangle$  in  $S$ , new leaf(0)  
 if  $t=1$  for all  $\langle x, t \rangle$  in  $S$ , new leaf(1)  
 Else  
 Best Feature, Feature set = Choose Best Feature( $S$ , Feature set,  $m'$ )  
 if Best Feature == Null, return  
 $S_0$  = all  $\langle x, t \rangle$  in  $S$ , where  $X_{Best}=0$   
 $S_1$  = all  $\langle x, t \rangle$  in  $S$ , where  $X_{Best}=1$   
 Return New node( $S$ , Grow Rand Best Tree( $S_0$ , Feature set,  $m'$ ), Grow Rand Best Tree( $S_1$ , Feature set,  $m'$ ))

קל להכליל עבור מאפיינים קטגוריים ונומרים.

**גבול ההחלטה של עצים:** עץ מחלק את מרחב הקלט למלבנים (היפר קופסאות) לא חופפים. אם המאפיינים נומרים, ניתן להגיע להפרדה מושלמת.

**4. Support Vector Machines =** קלסיפייר בינארי שאינו חייב להיות לינארי המבוסס על טריק kernel.

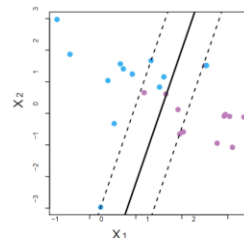
בהינתן קבוצת אימון  $D$  של דוגמאות נרצה ליצור מישור מפריד בעל השוליים המקסימליים. כאשר ממקסמים את השוליים, רק היפר-מישור אחד יכול להיות עם שוליים מקסימליים. אם נקודות האימון מחוץ לשוליים יזווגו, המישור המפריד לא יושפע אם נקודות האימון התומכות יזווגו, המישור המפריד ישתנה. נדמה את מישור ההפרדה (כולל השוליים) כשרוול / צינור ממימד  $n$ . נחפש soft margin classifier שאינו מכריח כל דוגמה להיות מעבר לשוליים וזאת על ידי כך שנאפשר שמיעוט

1. **One VS One** = בנה  $\frac{k(k-1)}{2}$  קלסיפיירים,

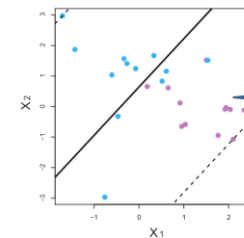
אחד לכל זוג קטגוריות, בהינתן דוגמא  $x$  הפעל את כל הקלסיפיי. פלט החיזוי: הקטגוריה בעלת רוב החיזויים

2. **One VS. Rest** = בנה  $k$  קלסיפיירים שמבדילים אם קלט  $x$  שייך לקטגוריה  $i$  או לכל יתר הקטגוריות. פלט החיזוי: הקלסיפיי שנותן את הביטחון המקסימלי לקטגוריה שלו:  $h_i(x)$  המקסימלי.

להלן ההשפעה של התקציב:



C קטן יאפשר רק margin צר, מספר קטן של נקודות תמך, וואריאנס גדול



C גדול מאפשר margin רחב; מספר גדול של נקודות תמיכה, וואריאנס קטן

גבולות החלטה לא לינארים, טריק kernel:

שיטת SVM משתמשת בטריק kernel שמאפשר הרחבת המרחב בו משתמש הקלסיפיי באופן שהוא בהרבה מיקרים יעיל יותר מבחינה חישובית וללא המרה בפועל של הווקטורים.

ההיפותזה:  $h(x) = \alpha_0 + \sum_{i=1}^m \alpha_i \langle x, x_i \rangle$  כלומר, מכפלה פנימית שהיא סוג של פונקציה דמיון בין שתי נקודות, לכן נוכל להחליף אותה בהכללה שלה:

$$h(x) = \alpha_0 + \sum_{i=1}^m \alpha_i K(x, x_i)$$

**גרעין kernel** = פונקציה המכמתת דמיון בין שני וקטורים. קיימים המון סוגי גרעינים, הנפוצים:

1. גרעין ליניארי, דמיון קורלטיבי

$$\text{Linear: } K(x, x_i) = x \cdot x_i$$

2. גרעין פולינומיאלי, גם בצורתו הפשוטה

ביותר  $d = 2$ , נותן הרבה כוח

$$\text{Polynomial: } K(x, x_i) = (x \cdot x_i)^d$$

3. גרעין גאוסיאני, נותן דמיון בעל צורת

פעמון. הווקטורים דומים אם מרחקם

האוקלידי קטן. ככל שהסיגמה קטנה,

הגאוסיאן צר, ואילו ככל שהסיגמה גדולה,

הגאוסיאן רחב

$$\text{Gaussian: } K(x, x_i) = e^{-\frac{1}{2\sigma^2} \|x - x_i\|^2}$$

דמיון בין אלגוריתם SVM לאלגוריתמים אחרים:

דומה לרגרסיה לוגיסטית עם רגולריזציה  $ridge$ , כשמשמשים בגרעין גאוסני מקבלים קירוב

לאלגוריתם  $KNN$ .

אלגוריתם SVM ליותר משתי קטגוריות:

האלגוריתם נבנה עבור קלסיפיקציה בינארית ולא

ניתן להרחבה בקלות ליותר מאשר 2 קטגוריות.

הדרכים בו ניתן להרחיב הן:

### למידה בלתי מונחת:

נציג שני אלגוריתמים של למידה שאיננה מונחת.

✓ **שיטת אשכולות Clustering =**

זוהי שיטה בלתי פרמטרית מיועדת עבור

קלסיפיקציה. בהינתן קבוצה  $D$  של דוגמאות אימון (ללא תיוגים) רוצים לחלק לקבוצות זרות שמכסות

את  $D$  ומכילות דוגמאות "דומות" כך שהנקודות שבתוך קלסטר תהינה דומות זו לזו ואילו נקודות בקלסטרים שונות תהינה שונות זו מזו.

**אלגוריתם k-means clustering:**

K-means Classifier ( $K, D, distance$ ):

- 1) add randomly each example in  $D$  to one of the  $K$  clusters
- 2) while there are change in clusters:
  - 2.1) calculate center  $c_i$  of each cluster
  - 2.2) for each example  $x$ , calculate the  $distance(x, c_i)$  and add  $x$  to the cluster that is closest

מרכז הקלסטר הוא ווקטור ממוצעי המאפיינים של הדוגמאות שבקלסטר. לשם ביצוע האלגוריתם:

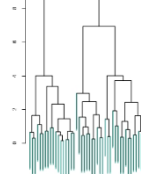
- 1) פונקצית מרחק לחישוב קרבה למרכז קלסטר, למשל מרחק אוקלידי.

2) מדד לשוני בתוך קלסטר **Within-Cluster-Variation**, למשל סכום הממוצעים של מרחקי

נקודה משאר הנקודות בקלסטר. חיסרון גדול של אלגוריתם זה הוא שיש לקבוע

מראש את מספר הקלסטרים. הפתרון לבעיה זו הוא שימוש בהיררכיה של קלסטרים.

**קלאסטרים היררכיים:**



מתחילים כשכל דוגמא בסט הנתונים היא קלסטר (עלה בעץ). בכל איטרציה מאחדים 2 קלסטרים שהם הכי קרובים עד שמקבלים בקלסטר אחד את כל הסט שוב. ברגע שינו גרף Dendograms נוכל ע"י חיתוך מלמעלה בגובה המתאים, לקבל כל מספר של קלסטרים שנרצה.

Hierarchical Classifier ( $K, D, ICD, WCV$ ):

- 1) add each example in  $D$  to a separate cluster
- 2) while more than one cluster left:
  - 2.1) calculate Inter-Cluster-Dissimilarity for each couple of clusters using ICD
  - 2.2) choose the two clusters with the minimal ICD and merge them into one cluster
  - 2.3) the height of the new united cluster is the result of the Within-Cluster-Variation using WCV

ישנן מספר דרכים לחשב אי דמיון בין קבוצות,

**linkage=inter-cluster-dissimilarity**

מקסימום אי דמיון: חשב את כל אי הדמיון עבור

כל זוג  $d(x, y)$  בין  $x \in C1$  ל- $y \in C2$  ותן כפלט את

השונויות המקסימלית.

**ממוצע:** חשב את כל אי הדמיון עבור כל זוג  $d(x, y)$

בין  $x \in C1$  ל- $y \in C2$  ותן כפלט את ממוצע השונויות.

**מינימום אי דמיון:** חשב את כל אי הדמיון עבור כל

זוג  $d(x, y)$  בין  $x \in C1$  ל- $y \in C2$  ותן כפלט את השונויות

המינימלית. נשים לב כי שיטה זו בעלת נטייה ליצור עץ לא מאוזן.

**מרכז:** תן כפלט את השונויות

$d(\text{center}(c1), \text{center}(c2))$  בין המרכזים של שני

האשכולות. בדרך כלל מרכז אשכול הוא הנקודה

הממוצעת. נשים לב כי שיטה זו בעלת נטייה לייצר

בעיות בויזואליזציה.

בשימוש באלגוריתמים אלו, יש צורך לנרמל את הנתונים, בשיטת  $sd$ . שיטת האשכולות לא תמיד עובדת- כאשר אין הפרדה טובה בין תת הקבוצות או כאשר נפח תת הקבוצות שונה משמעותית.

✓ **הפחתת מימדים PCA =**

בהינתן מדגם  $D$  של נתוני קלט  $n$  מימדי

$$x = (x_1, x_2, \dots, x_n)$$

לינאריות ששיעבירו את הקלט המקורי  $x$  ממימד  $n$

ל  $z$  ממימד  $m$  (קטן מ  $n$ ) תוך שימור של המידע

"החשוב". האלגוריתם מתבסס על הטלות ועל עקרון

מיקסום השונויות בו זמנית עם מזעור המרחק,

כלומר, המרחקים האוקלידיים של הנקודות מהציר מתמזערים באותו רגע בו מתמקסמת השונויות על

הציר.

הטרנספורמציות שמחפשים הם  $\phi_1, \phi_2, \dots, \phi_m$

כאשר  $z_i$  הוא ההטלה של דוגמא  $x$  על הציר ה- $i$ .

The score of the  $i$ th PC

$$z_i = \sum_{j=1}^n \phi_{j,i} x_j$$

$z = (z_1, z_2, \dots, z_m)$  הוא וקטור ההטלות של  $x$  על

מערכת צירים חדשה בה יש  $m$  צירים אורתוגונלים.

נרצה למצוא טרנספורמציות שישמרו מידע חשוב. כל

אחת מ  $m$  ה- $\phi_i$  הוא Principle Component.

הצירים החדשים הם הוקטורים העצמיים של

מטריצת ה- $Co$ -variance של  $D$ . הציר שנותן את

השונויות הכי גדולה בין ההטלות של הנקודות על

הציר הוא ה- $Principle$  Component הראשון שהוא

הוקטור העצמי עם הערך העצמי הדול ביותר של

מטריצה. הווקטור עם הערך העצמי הבא הכי גדול,

הוא הציר השני, כל הצירים מאונכים אילו לאילו.

יש לנרמל את הנתונים. התוצאות של האלגוריתם

יהיה תלויות בנרמול המאפיינים.



## דוגמאות: Naive Base

דוגמא:  
מטופל נבדק במעבדה והבדיקה חוזרת חיובית למחלת הסרטן. האם תבחר לעשות לו ניתוח מורכב כאשר יודעים את הרגישות והספציפיות של הבדיקה?

- ✓ recall: ידוע כי בדיקת מעבדה חיובית תתקבל ב 98% מהמיקרים בהם ישנה המחלה
- ✓ true negative rate: תוצאה שלילית נכונה תתקבל ב 97% מהמיקרים בהם המחלה איננה
- ✓ ידע מוקדם: 0.8% מהאוכלוסיה חולים במחלה

מנתונים אלה אנו מבינים:

$$\begin{aligned} P(-|cancer) &\approx 0.02, P(+|cancer) \approx 0.98 \\ P(-|\neg cancer) &\approx 0.97, P(+|\neg cancer) \approx 0.03 \\ P(cancer) &\approx 0.008, P(\neg cancer) \approx 0.992 \end{aligned}$$

ברור כי לסיווג ישנן שתי קטגוריות  $cancer$ ,  $\neg cancer$  וכן ישנו מאפיין אחד בוליאני שהוא תוצאת הבדיקה +, -.

כעת, נוכל לחשב את הסיווג על פי הקטגוריה בעלת ההסתברות המקסימלית מבין כל הקטגוריות האפשריות, כאמור, בשני אופנים:

1. לפי MAP:

$$\begin{aligned} y_{cancer} &= P(cancer|+) = P(+|cancer) \cdot P(cancer) = 0.98 \cdot 0.008 = 0.00784 \quad \text{I.} \\ y_{\neg cancer} &= P(\neg cancer|+) = P(+|\neg cancer) \cdot P(\neg cancer) = 0.03 \cdot 0.992 = 0.0298 \quad \text{II.} \end{aligned}$$

2. לפי ML:

$$\begin{aligned} y_{cancer} &= P(cancer|+) = P(+|cancer) = 0.98 \quad \text{I.} \\ y_{\neg cancer} &= P(\neg cancer|+) = P(+|\neg cancer) = 0.03 \quad \text{II.} \end{aligned}$$

ואז, נוכל לבחור היפותזה בהתאם לכל אחת מהגישות:

1. לפי MAP:  $0.0298 > 0.00784$  כלומר, הסיווג שיבחר הוא  $y_{\neg cancer}$ , לא חולה
  2. לפי ML:  $0.98 > 0.03$  כלומר, הסיווג שיבחר הוא  $y_{cancer}$ , חולה
- נשים לב כי דוגמה זו לא שיערה את ההסתברות עצמה, אלא רק נתנה סיווג שיבחר. במידה והיינו רוצים לחשב את השערוך להסתברות, היינו מחשבים:

1. עלינו לחשב את ההסתברות האחורית של לקבל בדיקה חיובית, ולכן:

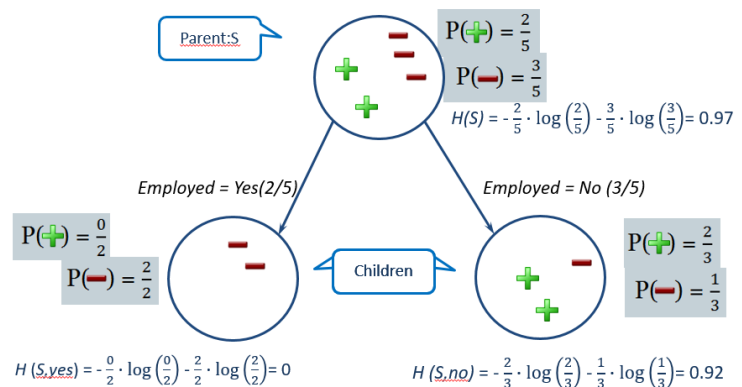
$$\begin{aligned} P(+ \cap cancer) &= p(+|cancer)p(cancer) \approx 0.98 \cdot 0.008 = 0.00784 \quad \text{I.} \\ P(+ \cap \neg cancer) &= p(+|\neg cancer)p(\neg cancer) \approx 0.03 \cdot 0.992 = 0.0298 \quad \text{II.} \\ P(+) &= P(cancer \cap +) + P(\neg cancer \cap +) \approx 0.00784 + 0.0298 = 0.03764 \quad \text{III.} \end{aligned}$$

2. כעת, על פי הנוסחה להתברות מותנית ניתן לחשב את השערוך:

$$\begin{aligned} P(cancer|+) &\approx \frac{P(+ \cap cancer)}{P(+)} = \frac{0.00784}{0.03764} \approx 0.21 \quad \text{I.} \\ P(\neg cancer|+) &\approx \frac{P(+ \cap \neg cancer)}{P(+)} = \frac{0.0298}{0.03764} \approx 0.79 \quad \text{II.} \end{aligned}$$

נשים לב כי החישוב הנ"ל מאמת את גישת MAP.

עצים:



נוכל לחשב את האנטרופיה המשוכללת:

$$P(D, E = yes) \cdot H(D, E = yes) + P(D, E = no) \cdot H(D, E = no) = \frac{2}{5} \cdot 0 + \frac{3}{5} \cdot 0.92 = 0.552$$

רווח המידע שיחושב הוא  $0.97 - 0.552 \approx 0.42$  כלומר, מאפיין "מועסק" הוסיף 0.42 כמות של מידע.

## רגרסיה פולינומיאלית:

$$\sum_{p=2}^P \binom{p+n-1}{p}$$

בהינתן  $m$  מאפיינים, להלן הנוסחה למספר המאפיינים שיש לחוסיף על מנת להעלות את דרגת הפולינום להיות  $p$  בעזרת  $C$ :

שערוך של ההסתברות הממוצעת לחיזוי נכון של קלסיפייר:

בהנחה שהקלסיפייר משתמש בפונקציית  $CE$  אז, שערוך ההסתברות הוא:

- חשב את השיגאה הממוצעת (מצע את כלל השיגאות על כל הדוגמאות), נסמן  $ce$
- חשב  $e^{-ce}$

## הוכחה הסתברותית לשימוש בנוסחת $CE$ :

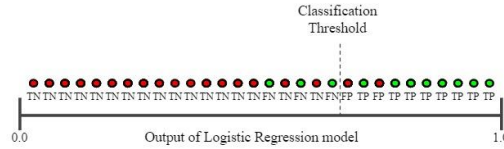
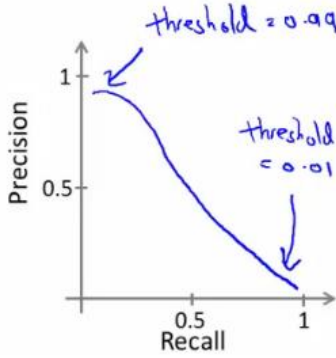
<p>נרצה למקסם את ה Likelihood שהיפותזה מסבירה את <math>D</math> (א.י). חוזה נכון את <math>t</math> של כל דוגמא ב <math>D</math>)</p> <p>ה-likelihood של <math>h_w</math> ביחס ל <math>D</math>, אינאלית <math>h_w</math> צריך לתת הסתברות 1 לדוגמאות חיוביות ו 0 לדוגמאות שליליות ואז שערוך ההסתברות הוא <math>l_D(w)=1</math> אם <math>h_w</math> אינו מדייק בחיזוי, <math>l_D(w) &lt; 1</math></p> <p>מיקסום ה likelihood מוצא את <math>w</math> שבו <math>hw</math> תותת הסתברויות גבוהות ל labels הנכונים,</p> $l_D(w) = \prod_{p:t_p=1} y_p \prod_{p:t_p=0} (1-y_p)$ <p><math>y_p = h_w(x_p)</math></p> <p>ממוצע גאומטרי לפי <math>m</math> דוגמאות יתן את ההסתברות הממוצעת לחזות נכון בדוגמא <math>x</math> בעזרת <math>w</math></p> <p>אימון: מחפשים <math>w</math> שימקסם את <math>l_D(w)</math></p> <p>ניתן למקסם את המכפלה באמצעות מיקסום <math>\log</math> שלה.</p> $w^* = \operatorname{argmax}_w \{\log(l_D(w))\}$ <p>לוג מכפלת ההסתברויות הוא סכום הלוגים</p> $\log \text{likelihood}_D(w) = \sum_p t_p \log(y) + \sum_p (1-t_p) \log(1-y)$	<p>הצדקה הסתברותית לשימוש ב <math>CE</math>:</p> <p>נרצה למצא היפותזה <math>h(w)</math> בעלת ההסתברות המקסימלית בהינתן <math>D</math>.</p> $\operatorname{argmax}_w \{P(h_w   D)\}$ <p>בהינתן <math>D</math>, מחפשים היפותזה <math>h_w</math> שמסבירה הכי טוב את <math>D</math></p> <p>א.י. שההסתברות <math>t_p = h_w(x_p)</math> דוגמא <math>p</math> ב <math>D</math>, היא מקסימלית מבין כל ה <math>h</math>ים האפשריים</p> <p>נראה כי מיעדור ה <math>CE</math> (תחת הנחות מסוימות) נותן את ההיפותזה הכי מסתברת</p>
<p>מדוע מיקסום loglikelihood עדיף?</p> <p>לוג מכפלת ההסתברויות הוא סכום הלוגים</p> $\log \text{likelihood}_D(w) = \sum_p t_p \log(y) + \sum_p (1-t_p) \log(1-y)$ <p>במקום למקסם מכפלה, ממקסמים את הלוג של המכפלה שהוא סכום של לוגים</p> <ul style="list-style-type: none"> <li>פעולה הסכום מהירה יותר ממכפלה</li> <li>בלוג יש פחות אובדן דיוק כאשר ההסתברויות קטנות מאוד.</li> <li>לעומת ההסתברויות, הלוגים הם מספרים גדולים בערכם המוחלט</li> <li>ההסתברות הממוצעת היא ממוצע גאומטרי (שורש <math>m</math>-י) אבל הפונקציה לממוצע רגיל של הלוגים</li> </ul>	<p>נרצה למצא היפותזה <math>h(w)</math> בעלת ההסתברות המקסימלית בהינתן <math>D</math>: <math>\operatorname{argmax}_w \{P(h_w   D)\}</math></p> <p>חוק: bayes</p> $p(h_w   D) = \frac{p(D h_w)P(h_w)}{P(D)}$ <p>נחשב את ה Likelihood: בהנחת אי תלות בין הדוגמאות:</p> <p>ההסתברות ש מודל <math>h_w</math> יחזוה נכון כל הדוגמאות ב <math>D</math>שווה למכפלת ההסתברות שכל דוגמא נחזית נכון</p> $l_D(h_w) = P(D h_w) = \prod_{p:D} p(t_p   x_p, h_w) = \prod_{p:t_p=1} p(t=1   x_p, h_w) \prod_{p:t_p=0} p(t=0   x_p, h_w)$ $= \prod_{p:t_p=1} h_w(x_p) \prod_{p:t_p=0} (1-h_w(x_p))$ <p>היפותזה מחשבת <math>t=1</math> הסתברות ל</p> $p(h_w   D) = \frac{P(h_w)}{P(D)} l_D(h_w)$
<p>Maximizing the LogLikelihood= Minimizing the CE</p> $\log \text{likelihood}_D(w) = \sum_p t_p \log(y_p) + \sum_p (1-t_p) \log(1-y_p)$ <p>אם ניקח ממוצע (של הדוגמאות ב <math>D</math>) ונקבל את ה <math>\log</math> של ההסתברות הממוצעת לחיזוי נכון של הדוגמאות ב <math>D</math></p> <p>נפחור את הסימן כדי לקבל מספרים חיוביים - שאותם נרצה למזער</p> <p>פונקציה זו נקראת Cross Entropy</p> $\text{CrossEntropy}(y, t) = -\frac{1}{m} \left( \sum_p t_p \log(y_p) + \sum_p (1-t_p) \log(1-y_p) \right)$ <p>מיעדור ה CrossEntropy היא דרך בלתי מאוד להתאים היפותזות <math>y=h_w(x)</math></p> <p>למטרת קלסיפיקציה</p> <p>ממצעים את ה <math>CE</math> כדי לקבל <math>-\log</math> ההסתברות הממוצעת</p> <p>מיעדור ה <math>CE</math> נקרא גם מיקסום של ה log-likelihood</p>	<p>נרצה למצא היפותזה <math>h(w)</math> בעלת ההסתברות המקסימלית בהינתן <math>D</math>: <math>\operatorname{argmax}_w \{P(h_w   D)\}</math></p> $p(h_w   D) = \frac{P(h_w)}{P(D)} l_D(h_w)$ $l_D(h_w) = \prod_{p:t_p=1} h_w(x_p) \prod_{p:t_p=0} (1-h_w(x_p))$ <p>מכיוון שרוצים למקסם את ההסתברות <math>p(h   D)</math>, ניתן להתעלם מה evidence - <math>p(D)</math> הקבוע שאינו תלוי ב <math>w</math></p> <p>נתעלם גם מ <math>P(h_w)</math> (הפריור של ההיפותזה): מחוסר ידע, מניחים שלכל היפותזות יש הסתברות שווה</p> $w^* = \operatorname{argmax}_w \{ \prod_{p:t_p=1} h_w(x_p) \prod_{p:t_p=0} (1-h_w(x_p)) \}$

הנחות בהוכחה: הקלטים לא תלויים סטטיסטית אחד בשני, מחוסר ידע מוקדם על היפותזות אפשריות הזנחנו את הפריור כלומר לכל ההיפותזות אותה הסתברות. אם ההנחות נכונות, ההיפותזה המתקבלת ממזעור  $CE$  היא ההיפותזה המסתברת ביותר בהינתן  $D$ .

## Precision – Recall Tradeoff

ככל שנקטין את הסף (נוזיז אותו שמאלה) :

1. Recall יעלה
2. Precision ירד



## שיטות רגישות לאי נרמול:

עבור שיטות אלו יש לנרמל תחילה את הנתונים :

1. PCA
2. KNN
3. Clustering
4. GD וכל מי שמשמש בו (רגרסיה לוגיסטית ולינארית)
5. SVM

במקרה של חריגים, עדיף לנרמל בעזרת std היות ואז הסטיית תקן תהיה 1.

<p><b>פירוק שגיאת MSE לניתנת להפחתה ולא:</b></p> $\begin{aligned} \text{MSE} &= E_{x,D}[(f(x) - h_D(x))^2] \\ &= E_{x,D}[(f(x) - \epsilon - h_D(x))^2] \\ &= E_{x,D}[(f(x) - h_D(x))^2 + 2(f(x) - h_D(x))\epsilon + (\epsilon - 0)^2] \\ &= E_{x,D}[(f(x) - h_D(x))^2] + 2E_{x,D}[(f(x) - h_D(x))\epsilon] + E_{x,D}[(\epsilon - 0)^2] \\ &= E_{x,D}[(f(x) - h_D(x))^2] + \text{Var}(\epsilon) = \text{ReducableError} + \text{IrreducibleError} \end{aligned}$	<p><b>סיבוכיות של GD:</b></p> <p>ריצה על כל ה-epochs: <math>O(e)</math></p> <p>ריצה על אוסף הדוגמאות: <math>O( D )</math></p> <p>ריצה על הפיצ'רים של דוגמה אחת: <math>O(d)</math></p> <p>ריצה על כל הפיצ'רים של כל דוגמאות האימון בכל ה-epochs: <math>O( D  * d * e)</math></p> <p>התשובה הנכונה: <math>O( D  * d * e)</math></p>	<p><b>סיבוכיות של המשוואה הנורמלית:</b></p> <p>הכפלה של <math>X^T</math> ב-<math>X</math>: <math>O(d^2 D )</math></p> <p>הכפלה של <math>X^T</math> ב-<math>t</math>: <math>O(d D )</math></p> <p>מציאת מטריצה הופכית ל-<math>X^T X</math>: <math>O(d^3)</math></p> <p>הכפלה של <math>(X^T X)^{-1}</math> ב-<math>X^T t</math>: <math>O(d^2)</math></p> <p>סה"כ: <math>O(d^2 D  + d D  + d^3)</math></p> <p>לאחר צמצום: <math>O(d^2 D  + d^3)</math></p> <p>התשובה הנכונה: <math>O(d^2 D  + d^3)</math></p>
<p><b>פירוק השגיאה הניתנת להפחתה לביאס ושונות:</b></p> $\begin{aligned} \text{ReducableError} &= E_{x,D}[(h_D(x) - f(x))^2] \\ &= E_{x,D}[(h_D(x) - \overline{h_D(x)} + \overline{h_D(x)} - f(x))^2] \\ &= E_{x,D}[(h_D(x) - \overline{h_D(x)})^2 + 2(h_D(x) - \overline{h_D(x)})(\overline{h_D(x)} - f(x)) + (\overline{h_D(x)} - f(x))^2] \\ &= E_{x,D}[(h_D(x) - \overline{h_D(x)})^2] + 2E_{x,D}[(h_D(x) - \overline{h_D(x)})(\overline{h_D(x)} - f(x))] + E_{x,D}[(\overline{h_D(x)} - f(x))^2] \\ &= \text{Variance}[h_D(x)] + \text{Bias}[h_D(x), f(x)]^2 \end{aligned}$	<p><b>נוסחת F Score שאיננה מאוזנת:</b></p> $F - \text{Measure} = \frac{1}{\alpha \cdot \frac{1}{\text{Precision}} + (1 - \alpha) \cdot \frac{1}{\text{Recall}}}$ <p>עבור <math>\alpha = 0.5</math> נקבל את <math>F1 \text{ score}</math></p>	<p><b>Feature Selection:</b></p> <p>במידה ולא משתמשים באלגוריתם חמדני, הסיבוכיות תהיה <math>O(2^n)</math></p> <p>כי נעבור על כל תתי הקבוצות האפשריות</p>

## רגולריזציה:

ככל שמעלים את קבוע ההענשה, השונות קטנה אילו הביאס גדל - ריבוע הביאס עולה בהתחלה לאט ואח"כ באופן מהיר יותר עד להתיצבותו. השונות יורדת בתמדה. ההכלה יורדת, מגיעה למינימום ועולה כשהביאס הופך לדומיננטי.

## אלגוריתם K-means:

אלגוריתם קליסטור של K-means מובטח שיתכנס כיוון שמספר הנקודות בדאטה סופי.