

Shelly Gamliel

TMI EX 3

Bayes Optimal and LDA

$$1) \forall x \in X \quad h_0(x) = \begin{cases} 1 & P(Y=1 | X=x) \geq \frac{1}{2} \\ -1 & \text{else} \end{cases}$$

$$h_0(x) = \operatorname{argmax}_{y \in \{-1, 1\}} P(X=x | Y=y) / P(Y=y) \quad : \text{B}$$

$P(Y=y | X=x)$ הינה פונקציית כיוון של x , y מושג y אם x מושג
レスポンס \rightarrow הינה פונקציית x , y מושג y אם x מושג

$l(y, \hat{y}) = 1_{\{y=\hat{y}\}}$ פונקציית נזק ביחס לטעות

$$= \begin{cases} 1 & \text{if } y = \hat{y} \\ 0 & \text{else} \end{cases}$$

$$\min_{y \in \{-1, 1\}} P(Y=y | X=x) \quad \text{הינה הינה פונקציית דבון הינה}$$

$$\max_{y \in \{-1, 1\}} P(Y=y | X=x) \quad -\Gamma \text{ פונקציית}$$

レスpond \rightarrow פונקציית $P(X=x | Y=y)$ פונקציית כיוון של y , x מושג y אם x מושג
פונקציית $P(Y=y)$ פונקציית כיוון של y , y מושג x מושג

$$P(Y=y | X=x) = P(X=x | Y=y) \frac{P(Y=y)}{P(X=x)} \quad \text{הנה שזאת חיקוי}$$

$$= \frac{P(X=x | Y=y) P(Y=y)}{\sum_{y \in Y} P(X=x | Y=y) P(Y=y)} = \frac{P(X=x | Y=y) P(Y=y)}{P(X=x | Y=1) P(Y=1) + P(X=x | Y=-1) P(Y=-1)}$$

פונקציית כיוון של y , x מושג y אם x מושג y אם x מושג
 $y \in \{-1, 1\}$

$$h_0(x) = 1 \iff \Pr(Y=1 | X=x) > \Pr(Y=-1 | X=x)$$

$$\iff \frac{\Pr(X=x | Y=1)}{\Pr(X=x)} \Pr(Y=1) > \frac{\Pr(X=x | Y=-1)}{\Pr(X=x)} \Pr(Y=-1)$$

$$\iff \Pr(X=x | Y=1) \Pr(Y=1) > \Pr(X=x | Y=-1) \Pr(Y=-1)$$

$$h_0(x) = -1 \iff \Pr(X=x | Y=-1) \Pr(Y=-1) > \Pr(X=x | Y=1) \Pr(Y=1)$$

ההypothesis $h_0(x)$ מוגדרת כפונקציית העדיפות של המודל. אם $\Pr(Y=1 | X=x) > \Pr(Y=-1 | X=x)$, אז $h_0(x) = 1$. אחרת, $h_0(x) = -1$.

$$h_0(x) = 1 \iff \Pr(Y=1 | X=x) > \frac{1}{2} \iff$$

$$\frac{\Pr(X=x | Y=1)}{\Pr(X=x)} \Pr(Y=1) \geq \frac{1}{2} \iff$$

$$\Pr(X=x | Y=1) \Pr(Y=1) \geq \Pr(X=x | Y=-1) \Pr(Y=-1)$$

$$\iff \Pr(X=x | Y=1) \Pr(Y=1) \geq \frac{2}{2} \Pr(X=x | Y=-1) \Pr(Y=-1)$$

$$\iff \Pr(X=x | Y=1) \Pr(Y=1) > \Pr(X=x | Y=-1) \Pr(Y=-1)$$

$$h_0(x) = -1 \iff \Pr(Y=-1 | X=x) < \frac{1}{2}$$

$$\iff \Pr(X=x | Y=1) \Pr(Y=1) < \Pr(X=x | Y=-1) \Pr(Y=-1)$$

$$\iff \Pr(X=x | Y=1) \Pr(Y=1) < \Pr(X=x | Y=-1) \Pr(Y=-1)$$

MAP estimation - classification problem

$$\min_{x \in X} P(X=x | Y=y) \quad \forall y \in Y$$

$$\max_{x \in X} P(X=x | Y=y) \quad \forall y \in Y$$

$$P(X=x | Y=y) = \frac{P(X=x)}{f_Y(y)} f_{Y|X=x}(y) \propto P(X=x) f_{Y|X=x}(y)$$

get y from x , x on \mathcal{X} \Rightarrow x is the class label of y

$$\max_{x \in X} P(X=x) f_{Y|X=x}(y)$$

binary classification using $X = \{x_1, x_2\}$ in

$$\iff \hat{x}^*(y) = x_1 \quad \text{if } P(X=x_1 | Y=y) > P(X=x_2 | Y=y)$$

$$P(X=x_1) f_{Y|X=x_1}(y) > P(X=x_2) f_{Y|X=x_2}(y)$$

$$2) P(Y=y | X=x) = \frac{f_{X|Y}(x|y) P(Y=y)}{P(X=x)}$$

nic popnug nivimim vici, $y \sim N(\mu, \Sigma)$ $\Rightarrow f_{X|Y}(x|y) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{1/2}} e^{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)}$

$$= C f_{X|Y}(x|y) P(Y=y) \quad \frac{1}{P(X=x)} = C \quad \text{noi}$$

$$= C \cdot P(Y=y) (2\pi)^{-\frac{d}{2}} |\det \Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

- נ $P(Y=y) \sim N(\mu, \Sigma)$ $\Rightarrow P(Y=y) = \frac{1}{\sqrt{2\pi^d |\Sigma|}} e^{-\frac{1}{2} (y-\mu)^T \Sigma^{-1} (y-\mu)}$

$$C' = (2\pi)^{-\frac{d}{2}} |\det \Sigma|^{-\frac{1}{2}} \cdot C$$

$$P(Y=y | X=x) = C' P(Y=y) \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

$$\log(P(Y=y | X=x)) = \underbrace{\log C'}_{y \in \{0, 1\}} + \log P(Y=y) - \frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)$$

$$\underset{y \in \{0, 1\}}{\operatorname{argmax}} \log P(Y=y) - \frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)$$

$$= \underset{y \in \{0, 1\}}{\operatorname{argmax}} \log P(Y=y) - \frac{1}{2} [x^T \Sigma^{-1} x + \mu^T \Sigma^{-1} \mu] + x^T \Sigma^{-1} \mu$$

$$= \underset{y \in \{0, 1\}}{\operatorname{argmax}} C' + \log(P(Y=y)) - \frac{1}{2} \mu^T \Sigma^{-1} \mu + x^T \Sigma^{-1} \mu$$

$$C' = -\frac{1}{2} x^T \Sigma^{-1} x \quad \text{noi}$$

$$= \underset{y \in \{0, 1\}}{\operatorname{argmax}} \ln P(Y=y) - \frac{1}{2} \mu^T \Sigma^{-1} \mu + x^T \Sigma^{-1} \mu$$

$$= \underset{y \in \{0, 1\}}{\operatorname{argmax}} \delta_y(x)$$

nic nivimim \Rightarrow $\delta_0(x) = 1$, $\delta_1(x) = 0$ \Rightarrow $\delta_0(x) = 1$, $\delta_1(x) = 0$

3) Write a formula to estimate μ_+ , μ_- and $P(y)$
based on training set $S = (x_1, y_1), \dots, (x_m, y_m)$

$$P(y=y) = \frac{1}{m} \sum_{i=1}^m y_i = \frac{y \text{ (sum of all } y_i \text{ in } S)}{m}$$

$$\mu_y = \frac{1}{\sum_{i=1}^m y_i} \sum_{i=1}^m x_i \quad \text{since } x_i \text{ have value } y \text{ if } y_i = y$$

$$\Sigma = \frac{1}{m-2} \sum_{y \in \{1, -1\}} \sum_{i:y_i=y} (x_i - \mu_y)^2 = \frac{1}{m-2} \sum_{y \in \{1, -1\}} \sum_{i:y_i=y} (x_i - \mu_y)(x_i - \mu_y)$$

we can calculate μ_+ and μ_- by summing up all x_i for $y_i = 1$ and $y_i = -1$ respectively.

4) Given some data find the mean and variance of spam

spam or not spam hidden info

spam \cap not spam \cap pic, spam \cap not pic, not spam \cap pic, not spam \cap not pic

False positive, False Negative : don't use if

classifiers do not make sense (e.g. $y = 0.5$)

true label \rightarrow \hat{y} = False Positive if $y = 1$ Type I error
 $\hat{y} = \{ \text{spam}, \text{not spam} \} / \{ y = \{\text{not spam}\} / \{ \text{spam} \}$

not spam - negative label, spam - positive label \leftarrow

spam \cap not spam label - $y = 1$ $\hat{y} = 1$ μ_0

spam \cap not spam label - $y = -1$ $\hat{y} = -1$

$y = 1$ $\hat{y} = 1$ true spam

$$y = -1, \hat{y} = 1$$

in box - $y = 1, \hat{y} = 1$: spam true

$$y = 1$$

$\hat{y} = 1$ TP

$$y = -1$$

FP (Type II error)

$$\hat{y} = -1$$

FN
Type II error

$$TN$$

SVM - Formulation

5) QP: $\underset{v \in \mathbb{R}^n}{\operatorname{argmin}} \frac{1}{2} v^T Q v + c^T v$

where $Q \in \mathbb{R}^{n \times n}$, $c \in \mathbb{R}^n$, $d \in \mathbb{R}^m$ $A \in \mathbb{R}^{m \times n}$

Write the Hard SVM problem as a QP problem in canonical form.

What are the values of Q, A, c, d that express this problem as QP in canonical form?

$$Q = 2 I_{n \times n} = \begin{pmatrix} 2 & & & \\ & 2 & & \\ & & \ddots & \\ & & & 2 \end{pmatrix}, w = (w_1, \dots, w_n)^T$$

$$V = (\bar{w}, b), c^T = 0, d = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$$

$$A = \begin{pmatrix} y_1 x_{11} & y_1 x_{12} & \dots & y_1 x_{1n} & y_1 \\ y_2 x_{21} & y_2 x_{22} & \dots & y_2 x_{2n} & y_2 \\ \vdots & \vdots & & \vdots & \vdots \\ y_m x_{m1} & y_m x_{m2} & \dots & y_m x_{mn} & y_m \end{pmatrix}$$

$$\underset{(w,b)}{\operatorname{argmin}} \|w\|^2 \quad \text{Hard SVM.} \quad \square$$

$$\text{st } \forall i \quad y_i (Kw, x_i) + b \geq 1$$

$$\frac{1}{2} V^T Q V + c^T V = \frac{1}{2} V^T I V + 0 \quad \Rightarrow \quad 0 \text{ P.O.}$$

$$= V^T V = w^T w = \|w\|^2$$

$$\underset{v \in \mathbb{R}^n}{\operatorname{argmin}} \frac{1}{2} v^T Q v + c^T v \Rightarrow \underset{(w,b)}{\operatorname{argmin}} \|w\|^2 \quad \text{min}$$

can be written as

$$Av < d \Rightarrow \begin{pmatrix} y_1 x_{11} & y_1 x_{12} & \dots & y_1 x_{1n} & y_1 \\ \vdots & \vdots & & \vdots & \vdots \\ y_m x_{m1} & y_m x_{m2} & \dots & y_m x_{mn} & y_m \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \\ b \end{pmatrix} \leq -1$$

$$-(y_1 x_{11} - y_1 x_{12} - \dots - y_1 x_{1n} y_1) \begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix} = -y_1 (x_1, w) + y_1 b \leq -1 \quad \text{for } i \leq n$$

$$-y_1 (x_1, w) + b \leq -1 \quad d = -1 \text{ P.O.}$$

$$\Rightarrow y_1 (x_1, w) + b \geq 1 \quad \text{Hard SVM!}$$

$$6) \text{ Soft SVM: } \underset{\omega, \xi}{\operatorname{argmin}} \frac{\lambda}{2} \|\omega\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i$$

$$\text{s.t. } \forall i \quad y_i \langle \omega, x_i \rangle \geq 1 - \xi_i$$

Show that this problem is equivalent to

$$\underset{\omega}{\operatorname{argmin}} \frac{\lambda}{2} \|\omega\|^2 + \frac{1}{m} \sum_{i=1}^m l^{\text{hinge}}(y_i \langle \omega, x_i \rangle)$$

where $l^{\text{hinge}}(\alpha) = \max\{0, 1 - \alpha\}$

$$l^{\text{hinge}}(y_i \langle \omega, x_i \rangle) = \max\{0, 1 - y_i \langle \omega, x_i \rangle\}$$

$$y_i \langle \omega, x_i \rangle \geq 1 - \xi_i \iff \xi_i \geq 0$$

$$1 - y_i \langle \omega, x_i \rangle \leq \xi_i$$

$$\forall i \quad y_i \langle \omega, x_i \rangle \geq 1 - \xi_i, \quad \xi_i \geq 0$$

$$\forall i \quad \xi_i = \max\{0, 1 - y_i \langle \omega, x_i \rangle\}$$

$$= h^{\text{hinge}}(y_i \langle \omega, x_i \rangle)$$

$$\xi_i \geq \max\{0, 1 - y_i \langle \omega, x_i \rangle\}$$

$$\xi_i > \max\{0, 1 - y_i \langle \omega, x_i \rangle\}$$

$$\xi'_i = \max\{0, 1 - y_i \langle \omega, x_i \rangle\}$$

$$\frac{\lambda}{2} \|\omega\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \iff \xi_i > \xi'_i$$

$$\geq \frac{\lambda}{2} \|\omega\|^2 + \frac{1}{m} \sum_{i=1}^m \xi'_i$$

מבחן קיון מינימום בפונקציית מילוי

$\xi_i \geq \xi'_i$ מוכיח כי ξ_i מינימום בפונקציית מילוי

$$\xi_i = \max\{0, 1 - y_i \langle \omega, x_i \rangle\}$$

$$\xi_i \geq \max\{0, 1 - y_i \langle \omega, x_i \rangle\}$$

$$\xi_i = \max\{0, 1 - y_i \langle \omega, x_i \rangle\}$$

$$\xi_i > \max\{0, 1 - y_i \langle \omega, x_i \rangle\}$$

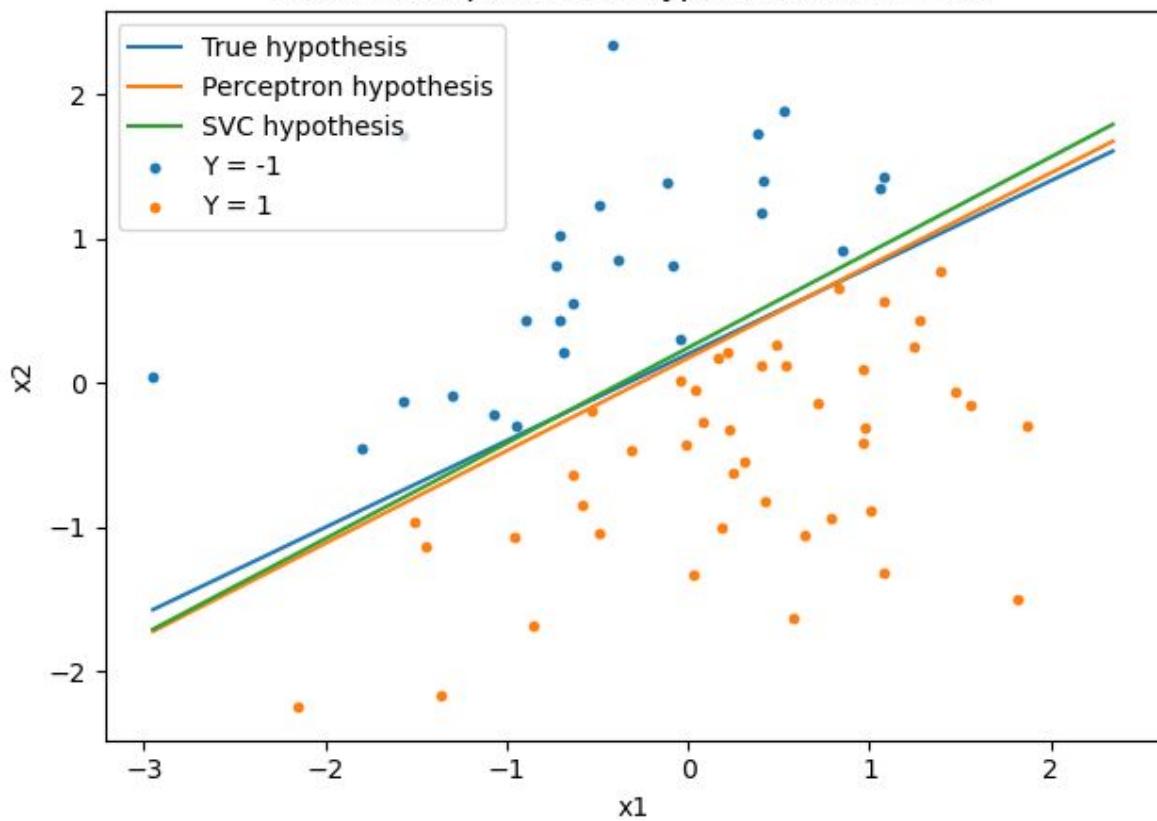
Iml Ex3

9. For each $m \in \{5, 10, 15, 25, 70\}$, draw m training points $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ and create the following figure:

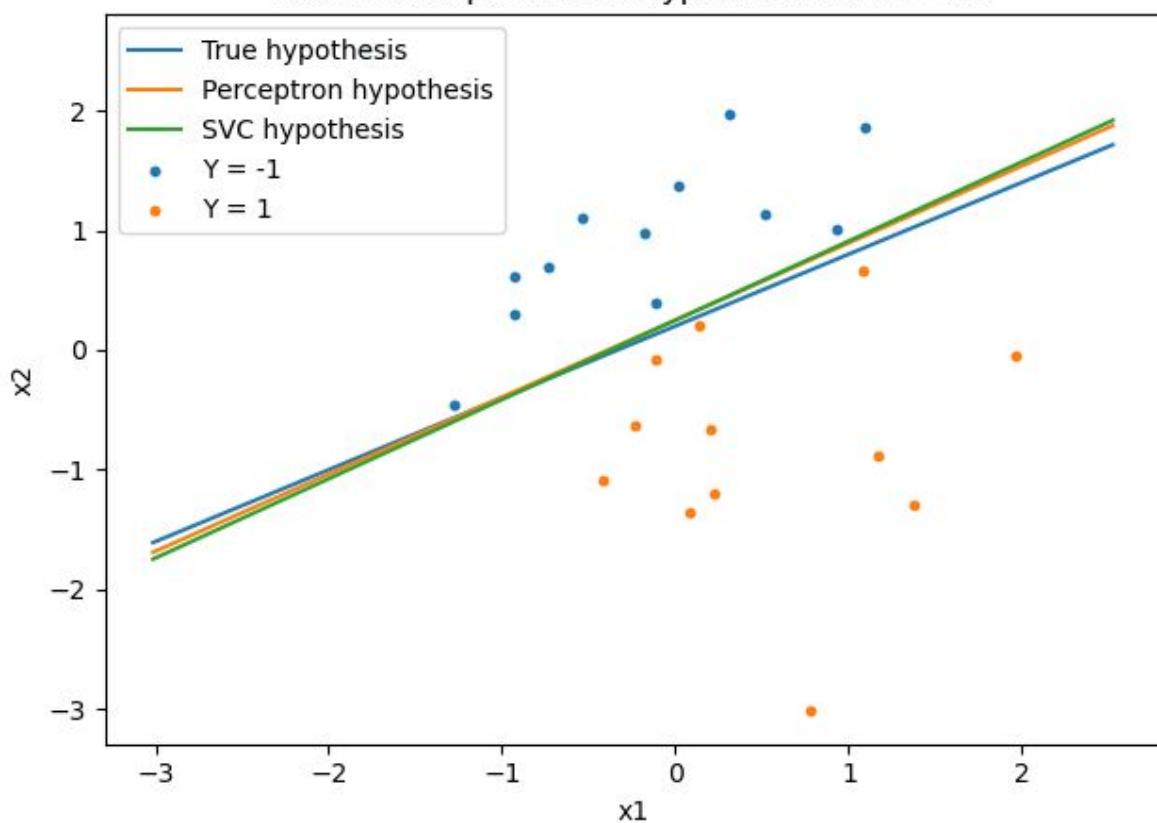
- The drawn data points, colored according to their labels (e.g., blue for positive labels and orange for negative ones)
- Add the hyperplane of the true hypothesis (the function f)
- Add the hyperplane of the hypothesis generated by the perceptron
- Add the hyperplane of the hypothesis generated by SVM
- Add a legend to explain which hyperplane is which

Add the 5 plots (one for each m) to your PDF file.

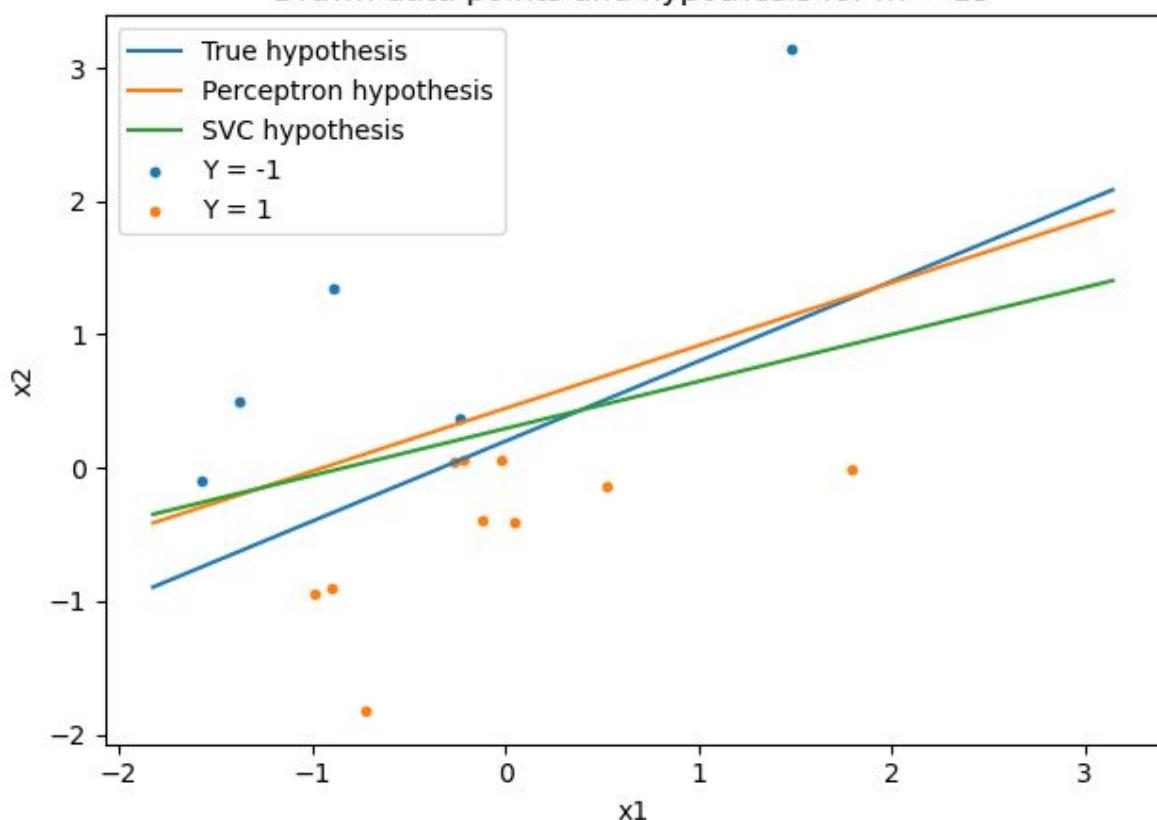
Drawn data points and hypothesis for $m = 70$



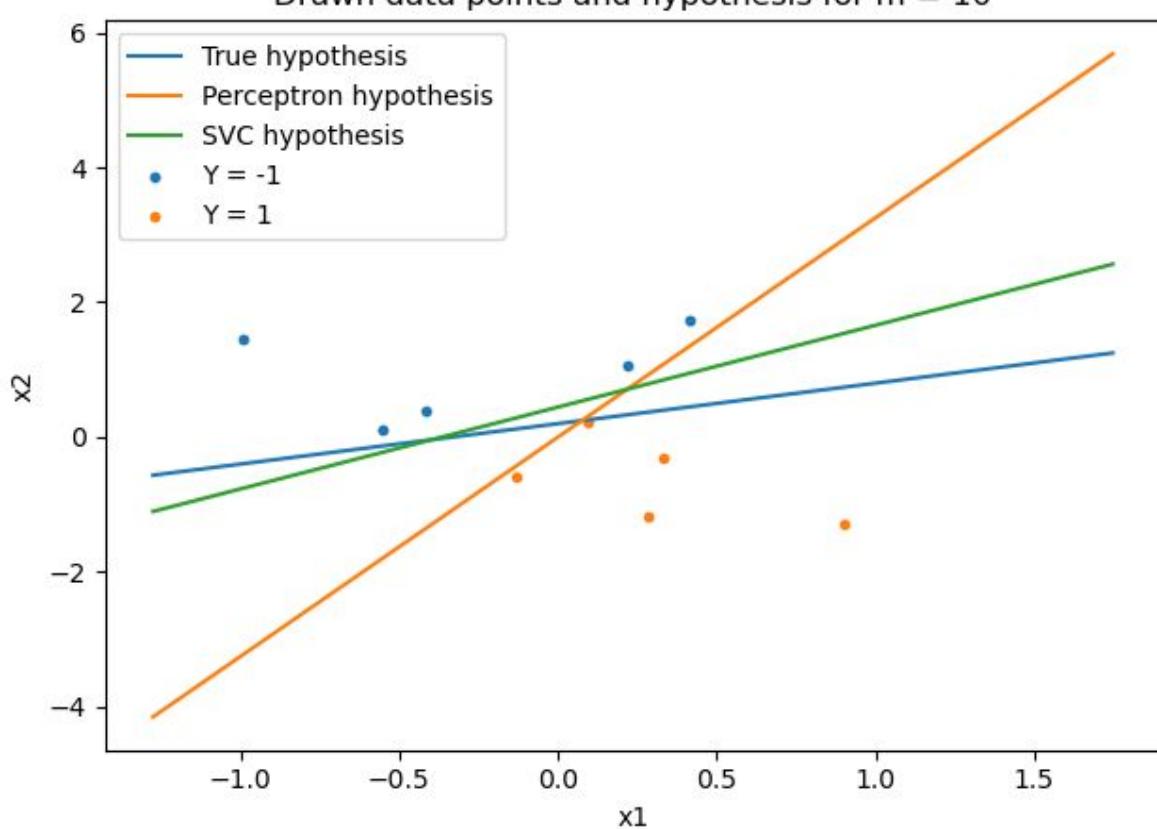
Drawn data points and hypothesis for $m = 25$



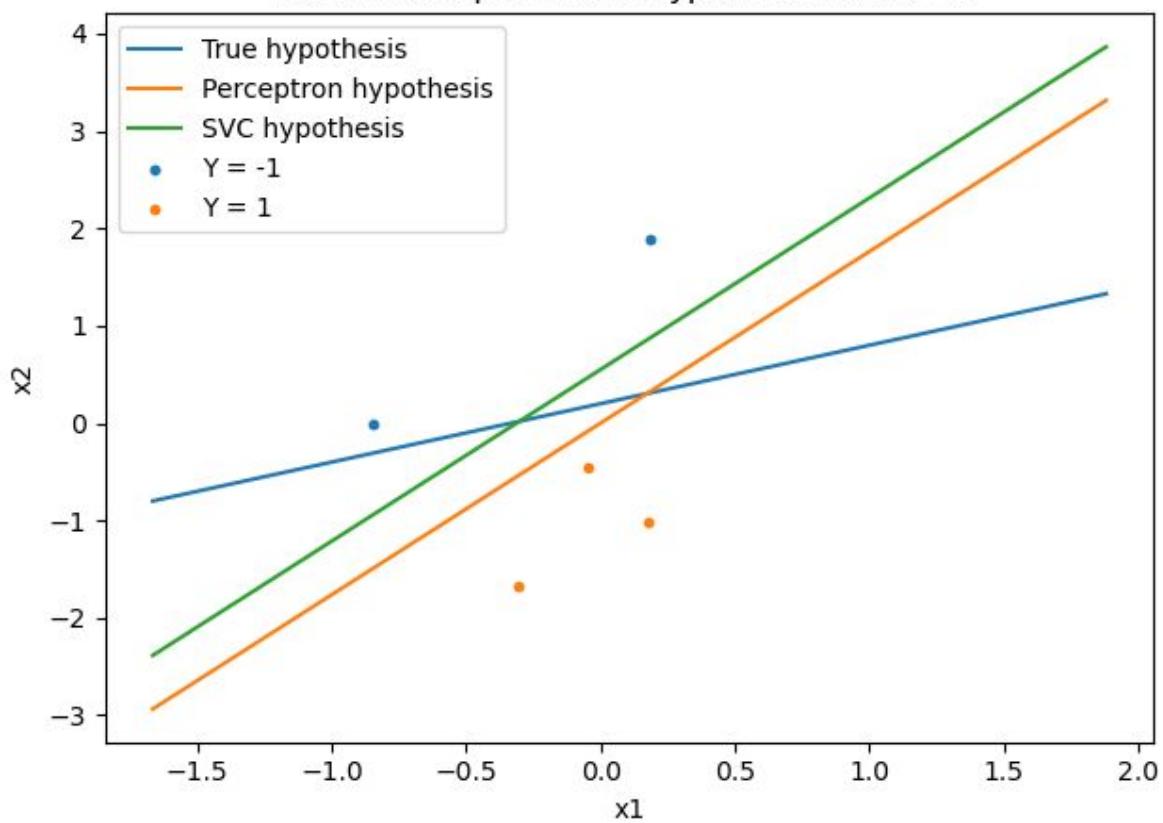
Drawn data points and hypothesis for $m = 15$



Drawn data points and hypothesis for $m = 10$



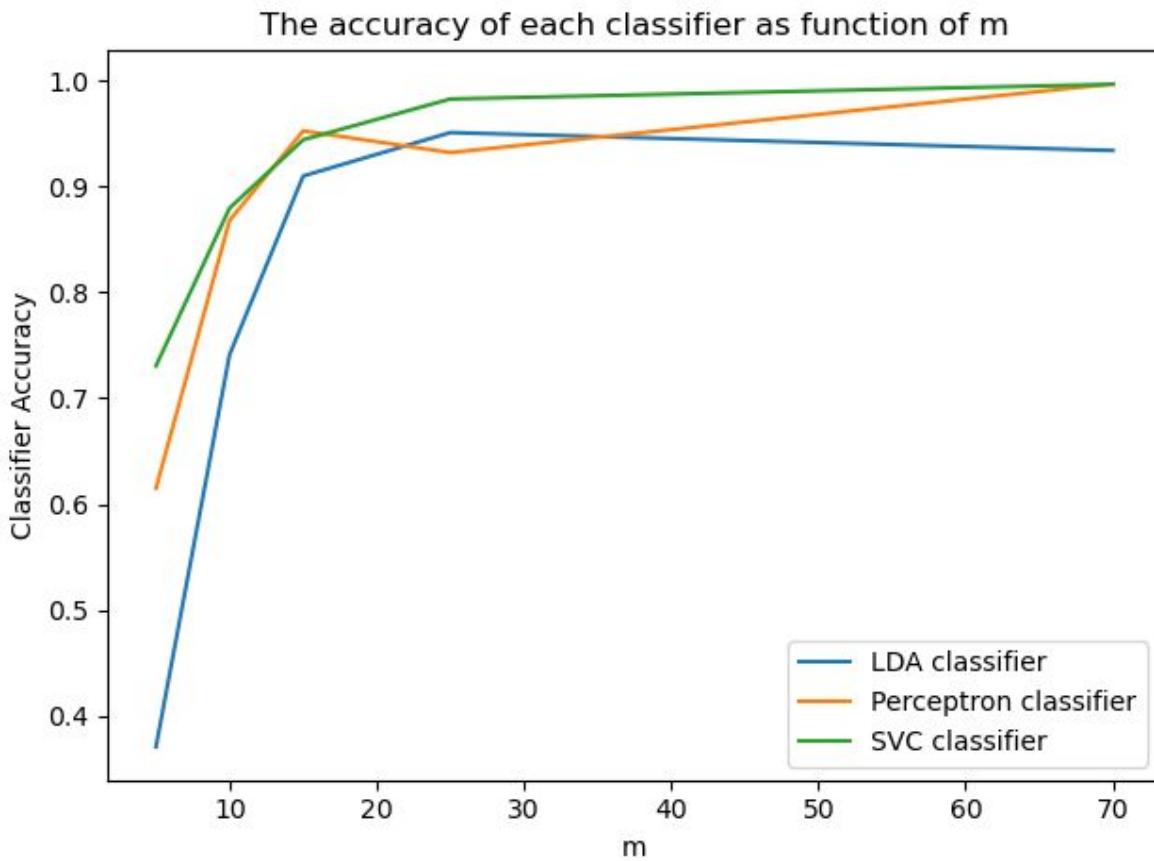
Drawn data points and hypothesis for m = 5



10. We'll now add a test set to compare three of the above algorithms. We'll use the following procedure:

- Draw training points $\{x_1, \dots, x_m\}$ from the distribution \mathcal{D} above and classify them according to a true hypothesis f (i.e. with labels y_1, \dots, y_m). Note: The training data should always have points from two classes. So if you draw a training set where no point has $y_i = 1$ or no point has $y_i = -1$ then just draw a new dataset instead, until you get points from both types.
- Draw k test points $\{z_1, \dots, z_k\}$ from the same distribution \mathcal{D} and calculate their true labels as well.
- Train a **Perceptron** classifier, an **SVM** classifier and an **LDA** classifier on $\{(x_i, y_i)\}_{i=1}^m$.
- Calculate the accuracy of these classifiers (the fraction of test points that is classified correctly) on $\{z_1, \dots, z_k\}$.

For each $m \in \{5, 10, 15, 25, 70\}$, repeat the above procedure 500 times with $k = 10000$ and save the accuracies (or just keep the mean accuracy, remember that the accuracy is a number between 0 to 1) of each classifier. Finally, plot the mean accuracy as function of m for each of the algorithms (SVM, Perceptron and LDA). Do not forget to add a legend to your graph. Add the plot to your PDF file.



11. Which classifier did better? why do you think that happened? No need for a formal argument, just explain what are the properties of the classifiers that cause these results.

המסוג מסווג וweis נתן את התוצאות הטובות ביותר ביותר מבחןית דיוון.
כאשר באופן כללי ה-lda classifier נותן תוצאות פחות טובות.
נשים לב כי עברו 70 דוגמאות ה-perceptron classifier וה-svm classifier מגיעים לדיוון גובה מאד.
אר עברו 25 דוגמאות ה-svm הגיע לדיוון טוב יותר.
נראה כי עברו מספר מספיק גדול של דוגמאות התוצאות של שני המסוגים קרובות, עברו מספר קטן יותר
ההווים מראה יכולות טובות יותר.

הסביר לגבי הביצועים של svm לעומת lda:

LDA: Assumes: data is Normally distributed. All groups are identically distributed. LDA is the best discriminator available in case all assumptions are actually met. QDA, by the way, is a non-linear classifier.

SVM: Generalizes the Optimally Separating Hyperplane. This hyperplane assumes that all groups are totally separable. SVM makes no assumptions about the data at all, meaning it is a very flexible method.

SVM classification is an optimization problem, LDA has an analytical solution. The optimization problem for the **SVM has a dual and a primal formulation that allows the user to optimize over either the number of data points** or the number of variables, depending on which method is the most computationally feasible.

LDA makes use of the entire data set to estimate covariance matrices and thus is somewhat prone to outliers. SVM is optimized over a subset of the data, which is those data points that lie on the separating margin.

הסביר לגבי הביצועים של svm לעומת perceptron:

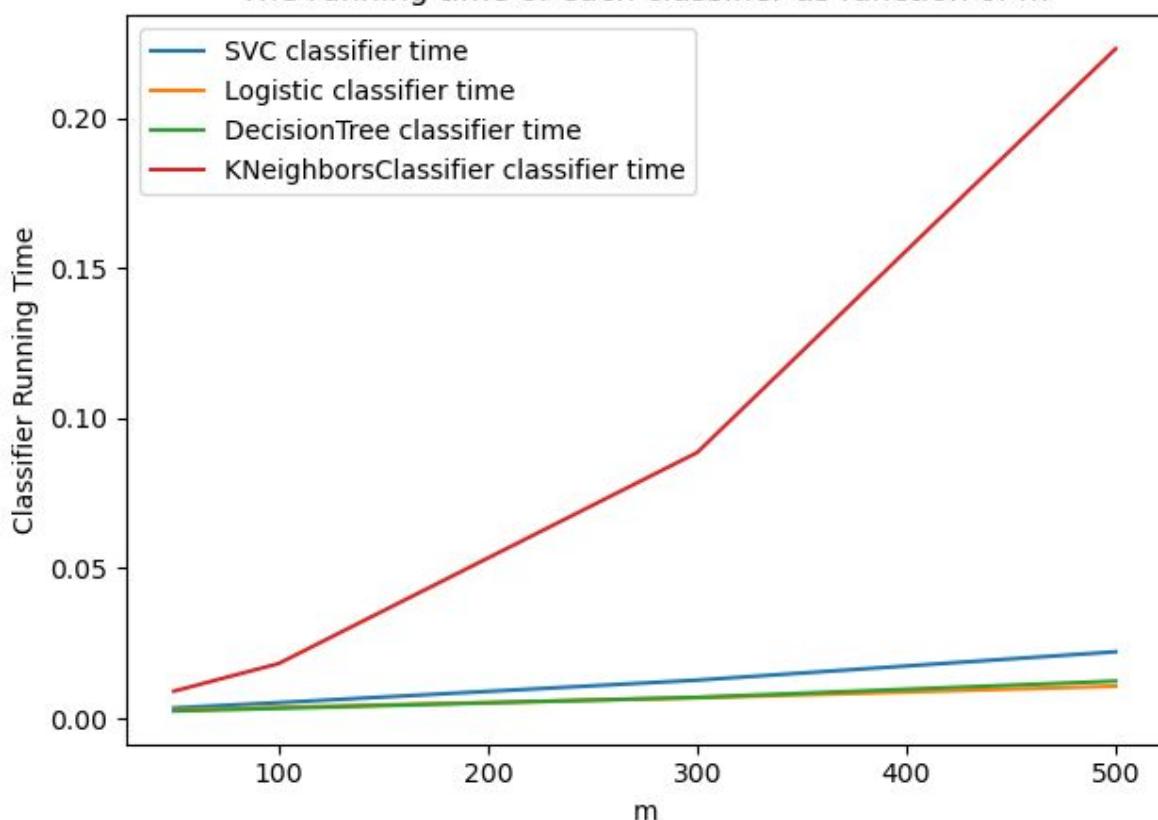
The major practical difference between a perceptron and SVM is that perceptrons can be trained online (i.e. their weights can be updated as new examples arrive one at a time) whereas SVMs cannot be. So, even though a SVM is usually a better classifier, perceptrons can still be useful because they are cheap and easy to re-train in a situation in which fresh training data is constantly arriving.

14. (similar to question 10, with different models) Consider the following procedure:

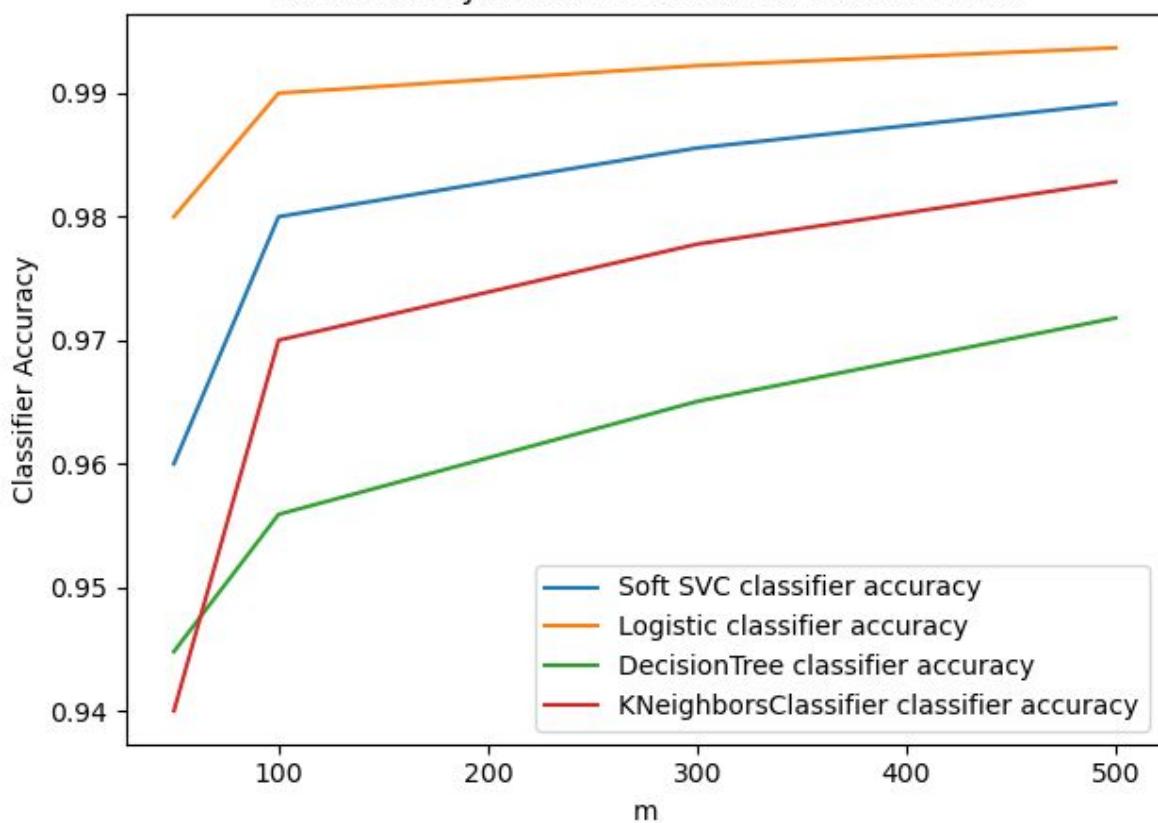
- Draw m training points $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ by choosing uniformly and at random from the train set. **Note:** The training data should always have points from two classes. So if you draw a training set where no point has $y_i = 1$ or no point has $y_i = -1$ then just draw a new dataset instead, until you get points from both types.
- Train a **logistic regression** classifier, an **Soft-SVM** classifier, a **decision tree** and a **k -nearest neighbors** classifier on $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$. You may choose any existing implementation of these algorithms as long as you note clearly in your solution which implementation you used. (Choose the regularization parameter for Soft-SVM and k for nearest neighbors to the best of your judgment - you're not expected to make the best possible choice at this stage. Although you are certainly encouraged to do your best and use side computations if you want.)
- Calculate their accuracy (the fraction of test points that is classified correctly) on the entire test set.

For each $m \in \{50, 100, 300, 500\}$, repeat the above procedure 50 times and save the elapsed running time and the accuracy (or just keep the mean accuracy, remember that the accuracy is a number between 0 to 1) of each classifier. Finally, plot the mean accuracy as function of m for each of the algorithms (SVM, Logistic regression, decision tree and nearest neighbors). Do not forget to add a legend to your graph. Add the plot to your PDF file. Discuss what you've seen about the difference in running time.

The running time of each classifier as function of m



The accuracy of each classifier as function of m



בgraf בו אנו מציגים את הדיק של כל אחד מהמודלים כפונקציה של מספר הדוגמאות מכל להבחן כי:
ה - logistic classifier נותן את הדיק הטוב ביותר (לא השפעה של מספר הדוגמאות) ביחס לשאר המודלים.

אחריו ה - soft svm classifier מוצג מוצמצם של 0.01 לכל היתר החל מ-100 דוגמאות ומעבר 0.020 דוגמאות.

מעבר 50 דוגמאות decision tree classifier נותן דיק טוב יותר לעומת k-nearest neighbors, אך החל מ 100 דוגמאות המצביע מותהף - k-nearest neighbors.

בgraf בו אנו מציגים את זמן הריצה נבחן כי kkk הינו בעל זמן הריצה הגבוה ביותר ביחס לאחרים וכך:

כל שמספר הדוגמאות גדול כך זמן הריצה גדול בצורה משמעותית.

מעבר המודלים האחרים כמו הדוגמאות משפיעות פחות על זמן הריצה לעומת kkn.
נבחן כי בהנתן דוגמה נמצא את k השכנים הקרובים ביותר ונחיזר תיג עלי הצביע הרוב.
כך שכל שיש יותר דוגמאות הבדיקה של השכנים הקרובים תגוזל זמן רב יותר.

השוואה:

Logistic regression vs SVM :

- SVM can handle non-linear solutions whereas logistic regression can only handle linear solutions.
- Linear SVM handles outliers better, as it derives maximum margin solution.
- Hinge loss in SVM outperforms log loss in LR.

Logistic Regression vs Decision Tree :

- Decision tree handles collinearity better than LR.
- Decision trees cannot derive the significance of features, but LR can.
- Decision trees are better for categorical values than LR.

Logistic Regression vs KNN :

- KNN is a non-parametric model, where LR is a parametric model.
- KNN is comparatively slower than Logistic Regression.
- KNN supports non-linear solutions where LR supports only linear solutions.
- LR can derive confidence level (about its prediction), whereas KNN can only output the labels.

K Nearest Neighbor (KNN) finds the “nearest examples” in the training data and chooses the label associated with the nearest examples. The algorithm takes a long time to compute. K will dominate by size of the class.

Logistic regression and support vector machines are closely linked. Both can be viewed as taking a probabilistic model and minimizing some cost associated with misclassification based on the likelihood ratio.

Logistic regression focuses on maximizing the probability of the data. The farther the data lies from the separating hyperplane (on the correct side), the happier LR is.

An SVM tries to find the separating hyperplane that maximizes the distance of the closest points to the margin (the support vectors).

Soft-SVM

- If training sample is not linearly separable, Hard-SVM returns no solution
- We no longer assume linearly separable. **Let's relax** the constraint to yield **soft-SVM**

$$\underset{\mathbf{w}, \xi}{\operatorname{argmin}} \left(\lambda \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \right)$$

s.t. $\forall i, y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 - \xi_i$ and $\xi_i \geq 0$

Summary - Soft SVM

- **Hypothesis class \mathcal{H} :** Half-spaces
- **Learning principle** (choosing $f \in \mathcal{H}$): Linear separator with Maximum Margin
- **Computational implementation:** Quadratic programming

Summary: Logistic regression

- **Hypothesis class:**

$$\mathcal{H}_{logistic}^d = \{\mathbf{x} \mapsto \pi(\langle \mathbf{x}, \mathbf{w} \rangle)\} .$$

- **Learning principle:** Maximum Likelihood
- **Computational implementation:** Convex optimization to maximize log-likelihood. Specialized iterative method / general convex solver
- So the maximum likelihood principle tells us to choose $h \in \mathcal{H}_{logistic}^d$, (equivalently, the vector $\mathbf{w} \in \mathbb{R}^{d+1}$) by finding

$$\hat{\mathbf{w}} := \underset{\mathbf{w} \in \mathbb{R}^{d+1}}{\operatorname{argmax}} \sum_{i=1}^m \left[y_i \langle \mathbf{x}_i, \mathbf{w} \rangle - \log \left(1 + e^{\langle \mathbf{x}_i, \mathbf{w} \rangle} \right) \right] .$$

LR gives calibrated probabilities that can be interpreted as confidence in a decision. LR gives us an unconstrained, smooth objective. LR can be used within Bayesian models. SVMs don't penalize examples for which the correct decision is made with sufficient confidence. This may be good for generalization.

Logistic Regression and trees differ in the way that they generate *decision boundaries* i.e. the lines that are drawn to separate different classes.

Decision Trees bisect the space into smaller and smaller regions, whereas Logistic Regression fits a single line to divide the space exactly into two. A single linear boundary can sometimes be limiting for Logistic Regression. If two classes are separated by a decidedly non-linear boundary, trees can better capture the division, leading to superior classification performance. However, when classes are not well-separated, trees are susceptible to overfitting the training data, so that Logistic Regression's simple linear boundary generalizes better.