



KTH Computer Science  
and Communication

## Exam in DD2421 Machine Learning 2023-03-16, kl 14.00 – 18.00

Aids allowed: *calculator, language dictionary.*

To take this exam you must be registered to this specific exam as well as to the course.

In order to pass this exam, your score  $x$  first needs to be 20 or more (out of 42, full point). In addition, given your points  $y$  from the Programming Challenge (out of 18, full point), the requirements on the total points,  $p = x + y$ , are preliminarily set for different grades as:

$$54 < p \leq 60 \rightarrow A$$

$$48 < p \leq 54 \rightarrow B$$

$$42 < p \leq 48 \rightarrow C$$

$$36 < p \leq 42 \rightarrow D$$

$$29 < p \leq 36 \rightarrow E \text{ (A pass is guaranteed with the required points for 'E'.)}$$

$$0 \leq p \leq 29 \rightarrow F$$

This exam consists of sections **A**, **B**, and **C**. **NB. Use different papers (answer sheets) for different sections.**

## A Graded problems

Potential inquiries to be addressed to Atsuto Maki.

### A-1 Terminology

(4p)

For each term (a–d) in the left list, find the explanation from the right list which *best* describes how the term is used in machine learning.

- |                               |  |
|-------------------------------|--|
|                               | 1) An approach to find useful dimension for classification                                 |
|                               | 2) Random strategy for area compression  |
| a) Dropout                    | 3) An approach to train artificial neural networks   |
| b) RANSAC                     | 4) Algorithm to learn with latent variables  |
| c) Occams's razor             | 5) A robust method to fit a model to data with outliers                                    |
| d) $k$ -fold cross validation | 6) A technique for margin maximization   |
|                               | 7) A technique for assessing a model by exploiting available data for training and testing |
|                               | 8) A principle to choose the simplest explanation  |

**Solution:** a-3, b-5, c-8, d-7

### A-2 Nearest Neighbor, Classification

(4p)

Suppose that we take a data set, divide it into training and test sets, and then try out two different classification procedures. We use *three-quarters* of the data for training, and the remaining *one-quarter* for testing. First we use *Adaboost* and get an error rate of 10% on the training data. We also get the average error rate (weighted average over both test and training data sets) of 12%. Next we use  $k$ -nearest neighbor (where  $k = 1$ ) and get an average error rate (weighted average over both test and training data sets) of 10%.

*Answer the following questions while motivating them.*

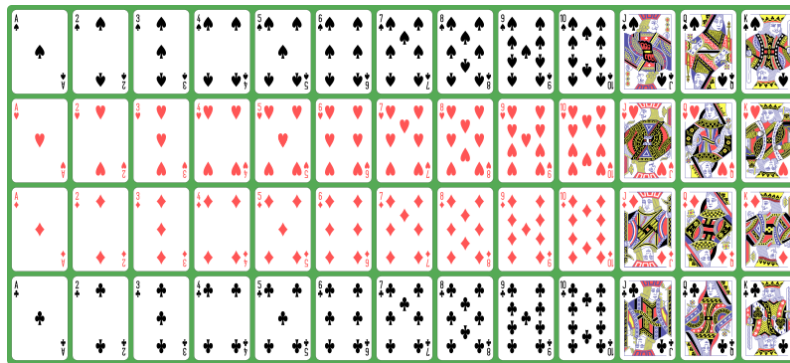
- a) What was the error rate with 1-nearest neighbor on the training set? (1p)
- b) What was the error rate with 1-nearest neighbor on the test set? (1p)
- c) What was the error rate with AdaBoost on the test set? (1p)
- d) Based on these results, indicate the method which we should prefer to use for classification of new observations, with a simple reasoning. (1p)

**Solution:**

- a) 0%. Training error for 1-NN is always zero.
- b) 40%. Given the answer in a), the testing error is 40%.
- c) 18%.
- d) AdaBoost, because it achieves lower error rate on the test data ( $18\% < 40\%$ ).

### A-3 Entropy and Decision Trees/Forests

(6p)



**Figure 1.** Playing cards consisting of 52 patterns.

- a) Consider decision trees, and indicate a correct one as the basic strategy for selecting a question (attribute) at each node.
- To minimize the expected reduction of the entropy.
  - To maximize the expected reduction of the entropy.
  - To minimize the expected reduction of gini impurity.

*Simply indicate your choice. (1p)*

- b) Briefly explain the two kinds of randomness involved in the design of Decision Forests? (2p)

Now imagine that you are playing with Cards and randomly sample *four cards* out of the pile of 52 cards (see Figure 1) *with replacement*, i.e. you sequentially draw a card but return it to the pile each time you have seen what it is.

- c) At each instance of drawing a card, what is the Shannon information content of the outcome with respect to the suit, one of {*Spades, Hearts, Diamonds, and Clubs*}, measured in bits? Motivate your answer. (1p)
- d) You play a game with a rule that you win if the suits of all the *four* cards are of *the same colour*, either *black or red*. Otherwise you lose. With respect to the outcome of this game, what is the expected information gain by drawing the first two card, i.e. by seeing (the suit colours of) the first card and the second card? Motivate your answer. (2p)

**Note:** if you do not have a calculator, answer with an expression but simplify it as much as possible.

**Solution:**

- a) ii
- b) Randomness (i) in generating bootstrap replicas and (ii) in possible features to use at each node.
- c) At each instance, it is  $\log_2 \frac{1}{1/4} = 2$  (bits).
- d) There are  $2^4 (= 16)$  patterns in terms of the combinations of the colours (with equal probability). For two of these you win (all black or all red), for the remaining cases you lose.

Let  $f(p_1, p_2) = -p_1 \log_2 p_1 - p_2 \log_2 p_2$

Considering how unpredictable the outcome of this game (win or lose) is, it can be measured in bits in terms of entropy:  $f(\frac{2}{16}, \frac{14}{16}) = -\frac{1}{8} \log_2 \frac{1}{8} - \frac{7}{8} \log_2 \frac{7}{8} \approx 0.54$

Two scenarios: the second card can bear the same colour as the first card, or different colour. These happen with probabilities,  $\frac{1}{2}$  and  $\frac{1}{2}$ , respectively. If the first two bear different colours, we know we will lose and hence the remaining entropy is zero. If the first two are of the same colour, we still have the chance for winning with the conditional probability of  $\frac{1}{4}$ , (the outcome unpredictable) with the entropy being  $f(\frac{1}{4}, \frac{3}{4}) \approx 0.81$ .

The information gain:  $0.54 - (\frac{1}{2} \cdot 0 + \frac{1}{2} \cdot 0.81) \approx 0.14$

#### A-4 Regression with regularization

(4p)

For a set of  $N$  training samples  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ , each consisting of input vector  $\mathbf{x}$  and output  $y$ , suppose we estimate the regression coefficients  $\mathbf{w}^\top (\in \mathbf{R}^d) = \{w_1, \dots, w_d\}$  in a linear regression model by minimizing

$$\sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \lambda \sum_{i=1}^d w_i^2$$

for a particular value of  $\lambda$ .

Now, let us consider different models trained with different values of  $\lambda$ , starting from a very large value (infinity) and *decreasing* it down to 0. Then, for parts a) through c), indicate which of i. through v. is most likely to be the case.

*Briefly justify each of your answers.*

- a) As we decrease  $\lambda$ , the variance of the model will:
  - i. Remain constant.
  - ii. Steadily increase.
  - iii. Steadily decrease.
  - iv. Decrease initially, and then eventually start increasing in a U shape.
  - v. Increase initially, and then eventually start decreasing in an inverted U shape. (1p)
- b) Repeat a) for the training error (residual sum of squares, RSS). (1p)
- c) Repeat a) for test RSS. (2p)

**Solution:** a)-ii, b)-iii, c)-iv

If  $\lambda$  was infinity, all  $w_i$  would be zero; the model is extremely simple, predicting a constant and has no variance with the prediction being far from actual value (thus with high bias). As we decrease  $\lambda$ , all  $w_i$  increase from zero toward their least square estimate values while steadily decreasing the *training* error to the ordinary least square RSS (and also decreasing bias) as the model continues to better fit training data. The values of  $w_i$  then become more dependent on training data, thus increasing the variance. The *test* error initially decreases as well, but eventually start increasing due to overfitting to the training data.

#### A-5 PCA, Subspace Methods

(4p)

We consider to solve a  $K$ -class classification problem with the Subspace Method and for that we compute a subspace  $\mathcal{L}^{(j)}$  ( $j = 1, \dots, K$ ) using training data for each class, respectively. That is, given a set of feature vectors (as training data) which belong to a specific class  $C$  (i.e. with an identical class label), we perform PCA on them and generated an orthonormal basis  $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$  which spans a  $p$ -dimensional subspace,  $\mathcal{L}$ , as the outcome.

Provide an answer to the following questions.

- a) Given that we compute  $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$  as eigenvectors of the auto-correlation matrix  $Q$  based on the training data, how should we choose the eigenvectors in relation to the corresponding eigenvalues of  $Q$ ? (1p)
- b) We have a new input vector  $\mathbf{x}$  whose class is unknown, and consider its projection length on each  $\mathcal{L}^{(j)}$ . Describe how the projection length is represented, using a simple formula. (2p)
- c) Given  $\mathbf{x}$ , we computed its projection length on each subspace as  $S^{(j)}$  ( $j = 1, \dots, K$ ), respectively. For a few classes among those with labels  $\{l, m, n\}$ , we had the following observations:  
 $S^{(l)}$  was the minimum of all  $S^{(j)}$ 's,  
 $S^{(m)}$  was the maximum of all  $S^{(j)}$ 's, and  
 $S^{(n)}$  was the closest to the average of all  $S^{(j)}$ 's.  
Which class should  $\mathbf{x}$  belong to? Simply choose a class label. (1p)

**Solution:**

- a) Choose the eigenvectors corresponding to the largest eigenvalues of  $Q$ .
- b)  $\sqrt{S}$  where  $S = \sum_{i=1}^p (\mathbf{x}, \mathbf{u}_i)^2$
- c)  $m$

## B Graded problems

Potential inquiries to be addressed to Bob Sturm.

### B-1 Identify Bayes' Theorem

(1p)

For two events  $A$  and  $B$ , which of the following is Bayes' Theorem:

- a)  $P(A|B) = P(A, B)P(A)/P(B)$
- b)  $P(A|B) = P(B|A)P(A, B)/P(B)$
- c)  $P(A|B) = P(B|A)P(A)/P(B)$
- d)  $P(A|B) = P(A, B)/P(B)$

**Solution:** (d)

### B-2 Conditioning

(1p)

For three events  $A$ ,  $B$ , and  $C$ , which of the following are not equivalent to  $P(A, B, C)$ :

- a)  $P(A, C|B)P(B)$
- b)  $P(A|B, C)P(B, C)$
- c)  $P(A, B|C)P(B)$
- d)  $P(A|C, B)P(B|C)P(C)$

**Solution:** (c)

### B-3 Shoogee picks

(2p)

Our dog Shoogee<sup>1</sup> (pictured) has a big box of  $N$  balls, all identical except for color: one is red and the others are brown. But since Shoogee is colorblind she can't see a difference. She likes to take the balls out one by one at random and lay them beside the box. Find an expression for the expected number of balls Shoogee takes out of the box before taking out the one red ball. *Show your work.*

**Solution:** Denote the expected number of balls taken given  $N$  balls as  $s(N)$ . This is defined

$$s(N) = \sum_{n=1}^N nP[A_n, B_n] = \sum_{n=1}^N nP[A_n|B_n]P[B_n] \quad (1)$$

where event  $A_n$  is ball  $n$  is red, and event  $B_n$  is  $n - 1$  previous balls are white, for  $1 < n \leq N$ . These distributions are given by:

$$P[A_n|B_n] = 1/(N - n + 1) \quad (2)$$

---

<sup>1</sup>About the name. Being non-native Swedish speakers, we were amused hearing people on the television saying "tjugo-noll-noll". Except they don't say "shoogo", but "shoogee". "Shoogo" sounds like a cleaning product. "Shoogee" sounds like a dog's name. So we named her "Shoogee Nulnul".



**Figure 2.** Shoogee (dog) with one of her favorite things in the world.

$$P[B_n] = \frac{\prod_{m=1}^{n-1} (N - m)}{N!} = \frac{1}{N(N - n)!} \quad (3)$$

Putting these together

$$s(N) = \frac{1}{N} \sum_{n=1}^N \frac{n}{(N - n + 1)(N - n)!} \quad (4)$$

#### **B-4 Inference by Minimizing Expected Loss**

(3p)

A good strategy for probabilistic inference given an observation  $x$  is to choose the value  $y \in \mathcal{Y}$  that minimizes the expected loss:

$$y^*(x) = \arg \min_{y' \in \mathcal{Y}} E_{Y|X=x}[\mathcal{L}(y, y')] = \arg \min_{y' \in \mathcal{Y}} \int_{y \in \mathcal{Y}} \mathcal{L}(y, y') Pr(y|X = x) dy \quad (5)$$

where  $\mathcal{L}(y, y')$  is a pre-defined loss function. Define  $\mathcal{L}(y, y') = (y - y')^2$  and assume  $Y|X = x \sim \text{Unif}[0, 2x]$ ,  $x > 0$ , that is a uniform distribution over the domain 0 and  $2x$ , for all positive  $x$ . Derive a formula for  $y^*(x)$  in this case. *Show your work.*

**Solution:** The conditional random variable  $Y|X = x$  is distributed

$$Pr[y|X = x] = \frac{1}{2x}, 0 \leq y \leq 2x \quad (6)$$

and zero otherwise. Plugging this and the loss into the equation above

$$y^*(x) = \arg \min_{y' \in [0, 2x]} \int_{y=0}^{2x} (y - y')^2 \frac{1}{2x} dy = \arg \min_{y' \in [0, 2x]} \int_{y=0}^{2x} (y - y')^2 dy \quad (7)$$

To evaluate the integral, define  $\Delta = y - y'$ , then  $d\Delta = dy$  and

$$\int_{\Delta=-y'}^{2x-y'} \Delta^2 d\Delta = \frac{1}{3} \Delta^3 \Big|_{\Delta=-y'}^{2x-y'} = \frac{1}{3} [(2x - y')^3 + y'^3] \quad (8)$$

Taking the partial derivative of this with respect to  $y'$  produces

$$\frac{\partial}{\partial y'} \frac{1}{3} [(2x - y')^3 + y'^3] = 3(2x - y')^2(-1) + 3y'^2. \quad (9)$$

Setting this to zero and solving for  $y'$  produces the result:

$$(2x - y')^2(-1) + y'^2 = -(4x^2 - 4xy' + y'^2) + y'^2 = -4x^2 + 4xy' = -x + y' = 0. \quad (10)$$

Finally,

$$y^*(x) = x. \quad (11)$$

## B-5 Regression

(1p)

Consider you have modelled  $\mathcal{D} = ((x, y)_i : x, y \in \mathbb{R})$  with

$$y = w_0 + w_1 x^2 + \epsilon \quad (12)$$

where  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ , and  $w_0 < 0$ . Graph a solution in the  $x$ - $y$  plane for  $w_1 < 0$  and a solution for  $w_1 > 0$ . Label the intersection on the  $y$ -axis.

**Solution:** Both graphs will feature a parabola with a  $y$ -offset at  $w_0$  (below the  $x$ -axis), but  $w_1 < 0$  will be open downward and  $w_1 > 0$  will be open upward.

## B-6 Probabilistic classification, part 1

(1p)

I have a dataset of labeled feature vectors,  $\mathcal{D} = ((\mathbf{x}, y)_i)$ , where  $y \in \{0, 1\}$  and  $\mathbf{x} = (x_1, x_2, x_3)$  is three dimensional. It looks like the first element of  $\mathbf{x}$  is one of three characters  $\{r, g, b\}$ ; the second element is one of two characters  $\{y, n\}$ , and the third element is a positive number. Which of the following is the Naive Bayes classifier for an observation  $\mathbf{x}$ ?

- a)  $y(\mathbf{x}) = \arg \max_{y \in \{0, 1\}} Pr(x_1, x_2, x_3 | y) Pr(y)$
- b)  $y(\mathbf{x}) = \arg \max_{y \in \{0, 1\}} Pr(x_1 | y) Pr(x_2 | y) Pr(x_3 | y) Pr(y)$
- c)  $y(\mathbf{x}) = \arg \max_{y \in \{0, 1\}} Pr(y | x_1) Pr(y | x_2) Pr(y | x_3) Pr(y)$
- d)  $y(\mathbf{x}) = \arg \max_{y \in \{0, 1\}} Pr(y | x_1, x_2, x_3) Pr(y)$

**Solution:** (b)

## B-7 Learning by Minimizing the Minimum Expected Loss

(3p)



A good strategy for probabilistic learning is to choose the parameters  $\theta$  of the conditional distribution  $Pr[y|X = x, \theta]$  that minimize the minimum expected loss:

$$\begin{aligned}\theta^* &= \arg \min_{\theta \in \Theta} \left[ \int_{x \in \mathcal{X}} \min_{y' \in \mathcal{Y}} E_{Y|X=x, \theta} [\mathcal{L}(y, y')] Pr(x) dx \right] \\ &= \arg \min_{\theta \in \Theta} \left( \int_{x \in \mathcal{X}} \left[ \min_{y' \in \mathcal{Y}} \int_{y \in \mathcal{Y}} \mathcal{L}(y, y') Pr(y|X = x, \theta) dy \right] Pr(x) dx \right) \quad (13)\end{aligned}$$

where  $\mathcal{L}(y, y')$  is a pre-defined loss function. Define  $\mathcal{L}(y, y') = (y - y')^2$  and assume  $Y|X = x \sim \text{Unif}[0, \theta x]$ ,  $\theta > 0$ , that is a uniform distribution over the domain 0 and  $\theta x$ , for all positive  $x$ . Show that in this case  $\theta^* \rightarrow 0$  is the best thing to do. *Show your work.*

**Solution:** The conditional random variable  $Y|X = x$  is distributed

$$Pr[y|X = x, \theta] = \frac{1}{\theta x}, 0 \leq y \leq \theta x \quad (14)$$

and zero otherwise. Let's first solve for the minimum expected loss:

$$\min_{y' \in \mathcal{Y}} \int_{y \in [0, \theta x]} (y - y')^2 \frac{1}{\theta x} dy = \frac{1}{\theta x} \min_{y' \in \mathcal{Y}} \int_{y=0}^{\theta x} (y - y')^2 dy \quad (15)$$

Define  $\Delta = y - y'$ , then  $d\Delta = dy$  and

$$\int_{\Delta=-y'}^{\theta x - y'} \Delta^2 d\Delta = \frac{1}{3} \Delta^3 \Big|_{\Delta=-y'}^{\theta x - y'} = \frac{1}{3} [(\theta x - y')^3 + y'^3] \quad (16)$$

Taking the partial derivative of this with respect to  $y'$  produces

$$\frac{\partial}{\partial y'} \frac{1}{3} [(\theta x - y')^3 + y'^3] = 3(\theta x - y')^2(-1) + 3y'^2. \quad (17)$$

Setting this to zero and solving for  $y'$  produces the result:

$$(\theta x - y')^2(-1) + y'^2 = -(\theta^2 x^2 - 2\theta x y' + y'^2) + y'^2 = -\theta^2 x^2 + 2\theta x y' = 0 \quad (18)$$

and so the minimum occurs when  $y' = \theta x/2$  for a given  $x$ . The minimum expected loss for a given  $x$  is thus

$$\frac{1}{\theta x} \left[ \frac{1}{3} [(\theta x - \theta x/2)^3 + (\theta x/2)^3] \right] = (\theta x)^2/12. \quad (19)$$

Now the problem reduces to finding the  $\theta$  that makes this small for all  $x$ . More formally, we need to solve

$$\theta^* = \arg \min_{\theta \in \Theta} \int_{x>0} (\theta x)^2 Pr(x) dx. \quad (20)$$

Clearly, we can make this smaller and smaller by moving  $\theta$  closer to zero, no matter how  $X$  is distributed – thus leading to a collapsed distribution for  $Y|X = x$ .

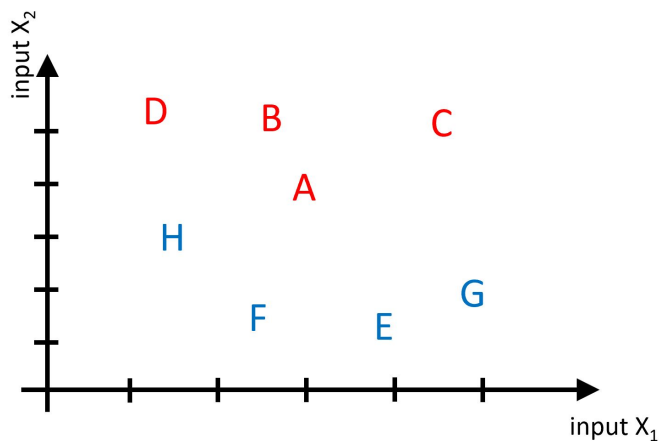
## C Graded problems

Potential inquiries to be addressed to Jörg Conradt.

### C-1 Support Vector Classification

(4p)

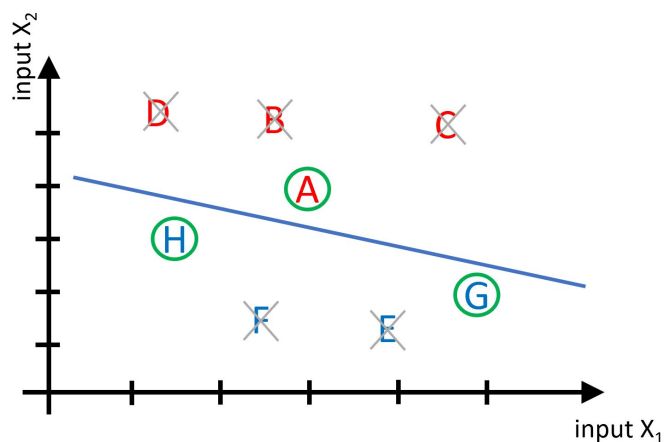
The following diagram shows a small data set consisting of four RED samples (A, B, C, D) and four BLUE samples (E, F, G, H). This data set can be linearly separated.



- a) We use a linear support vector machine (SVM) without kernel function to correctly separate the BLUE and the RED class. Which of the data points (A-H) can be removed for training without changing the resulting SVM decision boundary? Name the point(s) that can be removed, and justify with **KEYWORDS** (or very short sentences). (2p)
- b) Assume someone suggests using a non-linear kernel for the SVM classification of the above data set (A-H). Give one argument in favor and one argument against using non-linear SVM classification for such a data set. **USE KEYWORDS!** (2p)

**Solution:**

- a) The blue line shows the linear decision boundary between the two classes.



We can remove data points: **B, C, D, E, and F**.

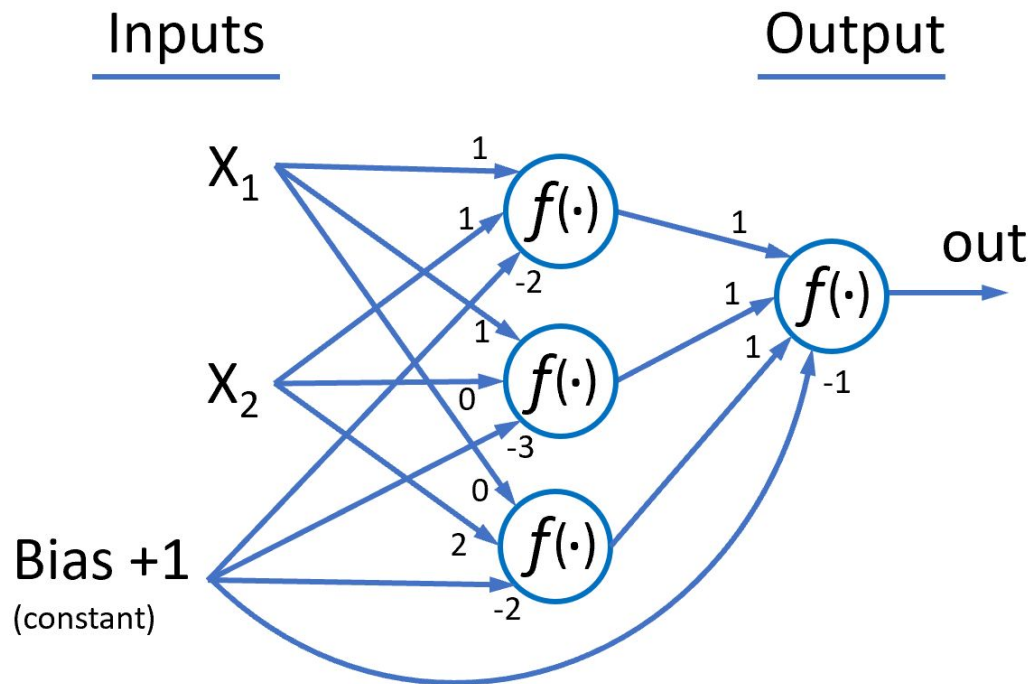
Justification: the other points (A, G, H) are closest to the decision boundary. If one of those points gets removed, the decision boundary changes. Therefore, those points are the required support vectors.

- b) at least one of each; max 1 point for (+) and 1 point for (-). Example arguments:
- + ) The decision boundary margin might get wider with a non-linear kernel.
  - + ) The same learning approach is likely to work for additional (possibly more complex) data.
  - ) More computing resources required.
  - ) The algorithm is more difficult to implement.

## C-2 Neuronal Networks

(4p)

The following diagram shows a simple neuronal network with step activation functions in all neurons. The weight for incoming signals are shown directly at the neurons.



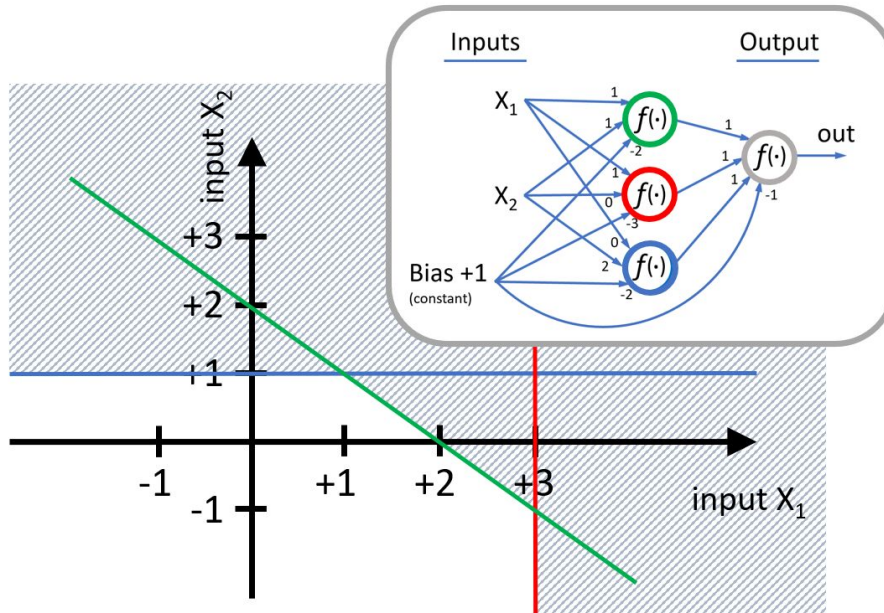
$f(\cdot)$  = step neuronal activation function

$$\text{step} \left( \sum_i (\text{inp}_i * w_i) \right) = \begin{cases} 1 & \text{if } (\cdot) \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

- a) Draw a 2D diagram of the input space (two dimensional plot of  $X_1$  vs.  $X_2$ ) and show for which area of the input space the network produces a positive output. (2p)
- b) Can this network be implemented in a single neuron with linear activation function (yes/no)? Explain in **KEYWORDS**. (1p)

- c) Assume all weights in the network double their value (the small numbers right next to the neuron input signal, including the weight for the constant bias input). What happens with the output? (1p)

**Solution:**



- a) Diagram shown above. Each separation line (colored) is given by one of the three neurons in the first layer. The output neuron is active if at least one of those neurons is active (OR function). Hence, the shaded region corresponds to the area of the input space ( $x_1, x_2$ ), where the network generates a positive output.
- b) The output function of this network is highly non-linear (a step with a piece-wise linear decision boundary in input space); hence a single linear neuron **CANNOT** implement this function.
- c) **No change at the output**, as all input (including bias) for all neurons doubles. Therefore, all total input that was below zero is still below zero; all total input that was above zero is still above zero. All neurons show the exact same output signals.