

FODA L29

---

Semester  
Review

~~Input~~  
~~Output~~

$(x, y)$   
supervised  
predict

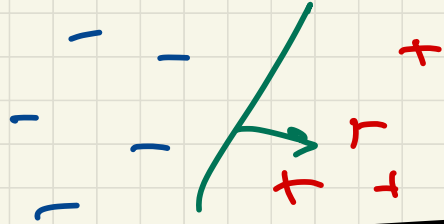
Can  
do  
cross  
validation

X

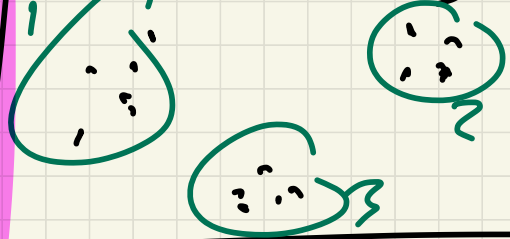
unsupervised  
find structure

{ - , + }  
"class"

Classification



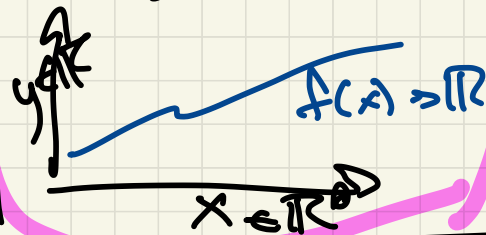
Clustering



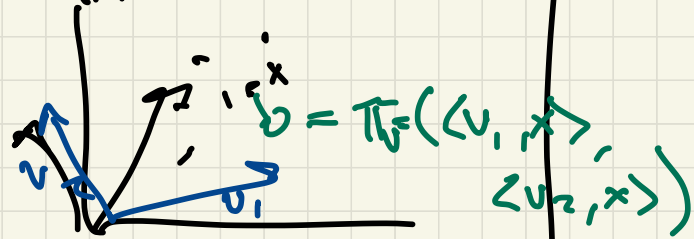
Real values

$\mathbb{R}$   
 $\mathbb{R}^n$

Regression



Dimensionality Reduction



1. Consider the random variables  $X$  and  $Y$  described by the joint probability table

	$X = 1$	$X = 2$	$X = 3$
$Y = 1$	0.10	0.05	0.10
$Y = 2$	0.30	0.25	0.20

$$\Pr[X=1 | Y=1] = \frac{0.1}{0.25} = 0.4$$

$$\Pr[Y=2] = 0.75$$

$$\Pr[X=3] = 0.3$$

$$\Pr[X=3 \cap Y=2] = 0.2 \neq (0.3)(0.75)$$

Derive the following values

(a)  $\Pr(X = 1) = 0.1 + 0.3 = 0.4$

(b)  $\Pr(X = 2 \cap Y = 1) = 0.05$

(c)  $\Pr(X = 3 | Y = 2) = \frac{0.2}{(0.3 + 0.25 + 0.2)} = \frac{0.2}{0.75}$

Compute the following probability distributions.

(d) What is the marginal distribution for  $X$ ?

(e) What is the conditional probability for  $Y$ , given that  $X = 2$ ?

$X=1$	$X=2$	$X=3$
0.4	0.3	0.3

$Y=1$	$0.05/0.3$
$Y=2$	$0.25/0.3$

Answer the following question about the joint distribution.

(f) Are random variables  $X$  and  $Y$  independent? **NO**

(g) Is  $\Pr(X = 1)$  independent of  $\Pr(Y = 1)$ ? **YES**

2. Consider two models  $M_1$  and  $M_2$ , where from prior knowledge we believe that  $\Pr(M_1) = 0.25$  and  $\Pr(M_2) = 0.75$ . We then observe a data set  $D$ . Given each model we assess the likelihood of seeing that data given the model as  $\Pr(D | M_1) = 0.5$  and  $\Pr(D | M_2) = 0.01$ . Now that we have the data, which model is has a higher probability of being correct?

Bayes' Rule  $\propto \Pr(D | M) \cdot P(M)$

$$\Pr(M | D) = \frac{\Pr(D | M) \cdot P(M)}{P(D)}$$

$$\Pr(M_1 | D) \propto \Pr(D | M_1) \cdot P(M_1) = (0.5) (0.25) = 0.125$$

$$\Pr(M_2 | D) \propto (0.01) (0.75) = 0.075$$

3. Assume I observe 3 data points  $x_1$ ,  $x_2$ , and  $x_3$  drawn iid from an unknown distribution. Given a model  $M$ , I can calculate the likelihood this each data point as  $\Pr(x_1 | M) = 0.5$ ,  $\Pr(x_2 | M) = 0.1$ , and  $\Pr(x_3 | M) = 0.2$ . What is the likelihood of seeing all of these data points, given the model  $M$ :  $\Pr(x_1, x_2, x_3 | M)$ ?

$$\begin{aligned} & \Pr(x_1, x_2, x_3 | M) && \text{independent} \\ &= \Pr(x_1 | M) \cdot \Pr(x_2 | M) \cdot \Pr(x_3 | M) \\ &= (0.5) (0.1) (0.2) = 0.001 \end{aligned}$$



5. Let  $X$  be a random variable that you know is in the range  $[-1, 2]$  and you know has expected value of  $E[X] = 0$ . Use the Markov Inequality to upper bound  $\Pr[X > 1.5]$ ?  
(Hint: you will need to use a change of variables.)

Markov  $X > 0$   $E[X] = \mu$

$$\Pr[X > \alpha] \leq \frac{E[X]}{\alpha}$$

$$Z = X + 1 \quad E[Z] = E[X] + 1 = 1$$

$$\Pr[X > 1.5] = \Pr[Z > 2.5] \leq \frac{E[Z]}{2.5} = \frac{1}{2.5}$$

6. Consider a matrix

$$A = \begin{bmatrix} 2 & 2 & 3 \\ -2 & 7 & 4 \\ -3 & -3 & -4 \\ -8 & 2 & 3 \end{bmatrix}$$

- (a) Add a column to  $A$  so that it is invertible. *full rank, square*
- (b) Remove a row from  $A$  so that it is invertible.
- (c) Is  $AA^T$  invertible? *no*
- (d) Is  $A^T A$  invertible? *yes*

$$AA^T = (4 \times 3)(3 \times 4) = 4 \times 4$$

*rank 3*

$$A^T A = (3 \times 4)(4 \times 3) = 3 \times 3$$

*rank 3*



7. Consider two vectors  $u = (0.5, 0.4, 0.4, 0.5, 0.1, 0.4, 0.1)$  and  $v = (-1, -2, 1, -2, 3, 1, -5)$ .

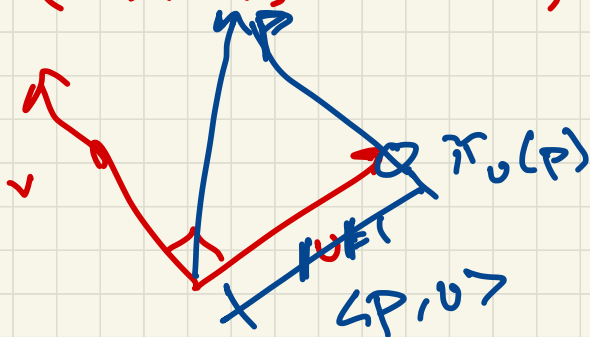
(a) Check if  $u$  or  $v$  is a unit vector. *not  $v$  yes  $||u|| \stackrel{?}{=} 1$*

(b) Calculate the dot product  $\langle u, v \rangle$ . *= 0*

(c) Are  $u$  and  $v$  orthogonal? *Yes*

$$\begin{aligned} & \sqrt{0.5^2 + 0.4^2 + 0.4^2 + 0.5^2 + 0.1^2 + 0.4^2 + 0.1^2} \\ & \quad 0.25 \quad 0.16 \quad 0.16 \quad 0.25 \quad 0.01 \quad 0.16 \quad 0.01 = 1 \end{aligned}$$

$$\begin{aligned} \langle u, v \rangle &= (-1) \cdot (0.5) + (-2) \cdot (0.4) + (1) \cdot (0.4) + \\ &= 0 \end{aligned}$$



8. Consider a matrix  $A \in \mathbb{R}^{n \times 4}$ . Each row represents a customer (there are  $n$  customers in the database). The first column is the age of the customer in years, the second column is the number of days since the customer entered the database, the third column is the total cost of all purchases ever by the customer in dollars, and the last column is the total profit in dollars generated by the customer.

For each of the following operations, decide if it is **reasonable** or **unreasonable**.

(a) Run simple linear regression using the first three columns to build a model to predict the fourth column. **Yes**

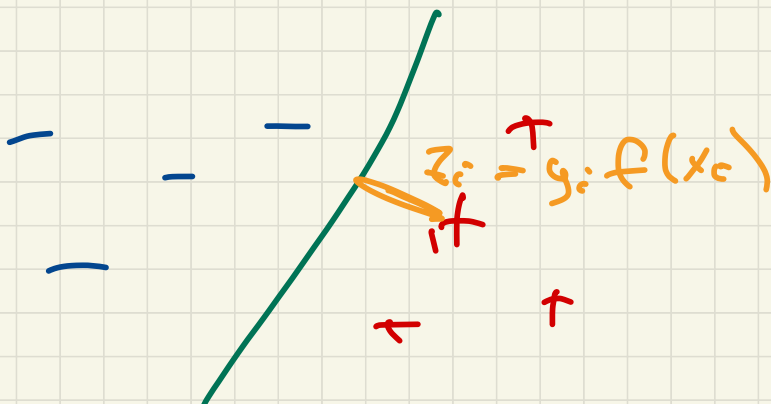
(b) Use  $k$ -means clustering to group the customers into 4 types using Euclidean distance between rows as the distance. **No**

dollars, years, days

(c) Use PCA to find the best 2-dimensional subspace, so we can draw the customers in a  $\mathbb{R}^2$  in way that has the least projection error. **No**

(d) Use the linear classification to build a model based on the first three columns to predict if the customer will make a profit +1 or not -1. **Yes**

supervised.



9. Consider a data set  $(X, y)$  where  $X \in \mathbb{R}^{n \times 3}$  we decompose into a test and a training data set  $(X_{\text{train}}, y_{\text{train}})$ . Assume that  $X_{\text{train}}$  is not just a subset of  $X$ , but also pads/prepends a column of all 1s. We build a linear model

$$\alpha = (X_{\text{train}}^T X_{\text{train}})^{-1} X_{\text{train}}^T y_{\text{train}}.$$

where  $\alpha \in \mathbb{R}^4$ . The remaining two testing data points are  $(x_1, y_1)$  and  $(x_2, y_2)$ , where  $x_1, x_2 \in \mathbb{R}^3$ . Explain (write a mathematical expression) to use this test data to estimate the generalization error. That is, if one new data point arrives  $x$ , how much squared error would we expect the model  $\alpha$  to have compared to the unknown true value  $y$ ?

$$\frac{(\langle \alpha, (1, x_1) \rangle - y_1)^2 + (\langle \alpha, (1, x_2) \rangle - y_2)^2}{2}$$

10. Consider a function  $f(x, y)$  with gradient  $\nabla f(x, y) = (x - 1, 2y + x)$ . Starting at a value  $(x = 1, y = 2)$ , and a learning rate of  $\gamma = 1$ , execute one step of gradient descent.

$$\begin{aligned}\nabla f(x=1, y=2) &= (1-1, 2(2)+1) \\ &= (0, 5)\end{aligned}$$

GD

$$(x', y') \leftarrow (x=1, y=2) - \gamma \nabla f(x=1, y=2)$$

$$\begin{aligned}& (1, 2) - (1) (0, 5) \\ &= (1-0, 2-5) \\ &= (1, -3)\end{aligned}$$

11. Consider running gradient descent with a fixed learning rate  $\gamma$ . For each of the following, we plot the function value over 10 steps (the function is different each time). Decide whether the learning rate is probably **too high**, **too low**, or **about right**.

- (a)  $f_1$ : 100, 99, 98, 97, 96, 95, 94, 93, 92, 91      too low
- (b)  $f_2$ : 100, 50, 75, 60, 65, 45, 75, 110, 90, 85      too high
- (c)  $f_3$ : 100, 80, 65, 50, 40, 35, 31, 29, 28, 27.5, 27.3      right
- (d)  $f_4$ : 100, 80, 60, 40, 20, 0, -20, -40, -60, -80, -100      too low

12. Consider a matrix  $A \in \mathbb{R}^{8 \times 4}$  with squared singular values  $\sigma_1^2 = 10$ ,  $\sigma_2^2 = 5$ ,  $\sigma_3^2 = 2$ , and  $\sigma_4^2 = 1$ .

- $v_4$
- (a) What is the rank of  $A$ ?  $4 = \# \text{ non-zero Sing. val.}$
- (b) What is  $\|A - A_2\|_F^2$ , where  $A_2$  is the best rank-2 approximation of  $A$ .  $= \sigma_3^2 + \sigma_4^2 = 2 + 1 = 3$
- (c) What is  $\|A - A_2\|_2^2$ , where  $A_2$  is the best rank-2 approximation of  $A$ .  $= \sigma_3^2 = 2$
- (d) What is  $\|A\|_2^2$ ?  $= 10$
- (e) What is  $\|A\|_F^2$ ?  $= 10 + 5 + 2 + 1 = 17$

Let  $v_1, v_2, v_3, v_4$  be the right singular vectors of  $A$ .

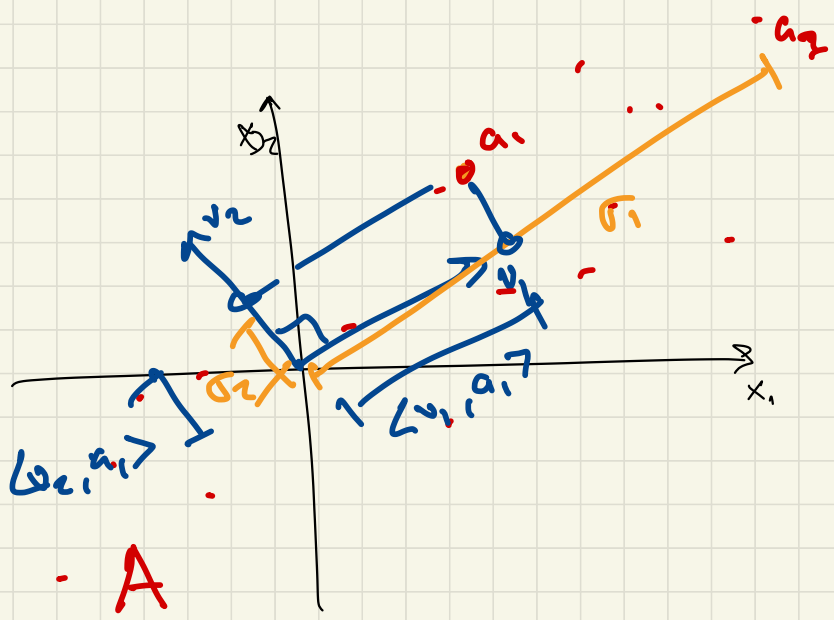
- (f) What is  $\|Av_2\|^2$ ?  $= \sigma_2^2 = 5$
- (g) What is  $\langle v_1, v_3 \rangle$ ?  $= 0$
- (h) What is  $\|v_4\|^2$ ?  $= 1$

Let  $a_1 \in \mathbb{R}^4$  be the first row of  $A$ .

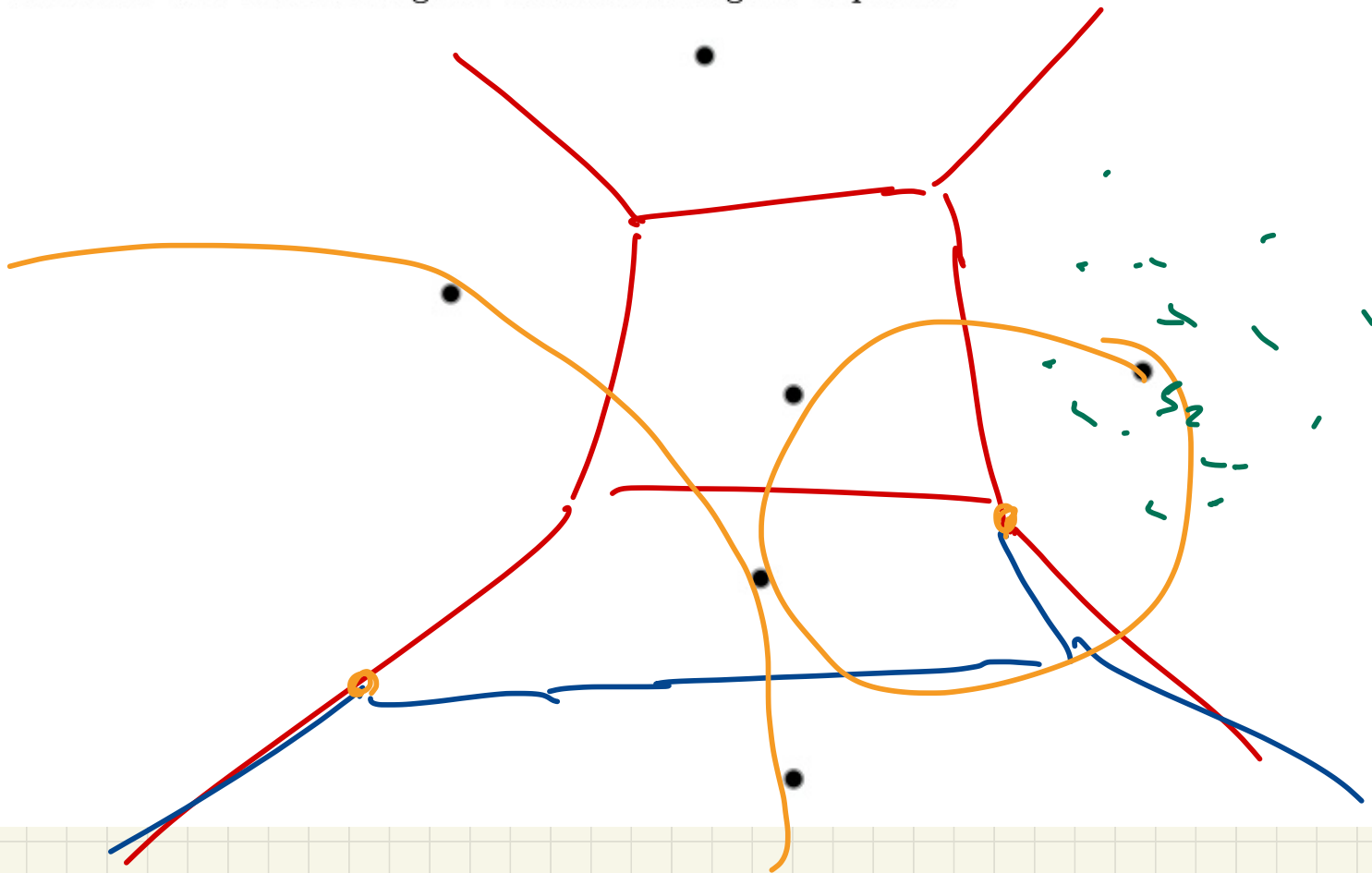
- (i) Write  $a_1$  in the basis defined by the right singular vectors of  $A$ .

$$\left( \langle v_1, a_1 \rangle, \langle v_2, a_1 \rangle, \langle v_3, a_1 \rangle, \langle v_4, a_1 \rangle \right) \in \mathbb{R}^4$$

$V^T a_1$

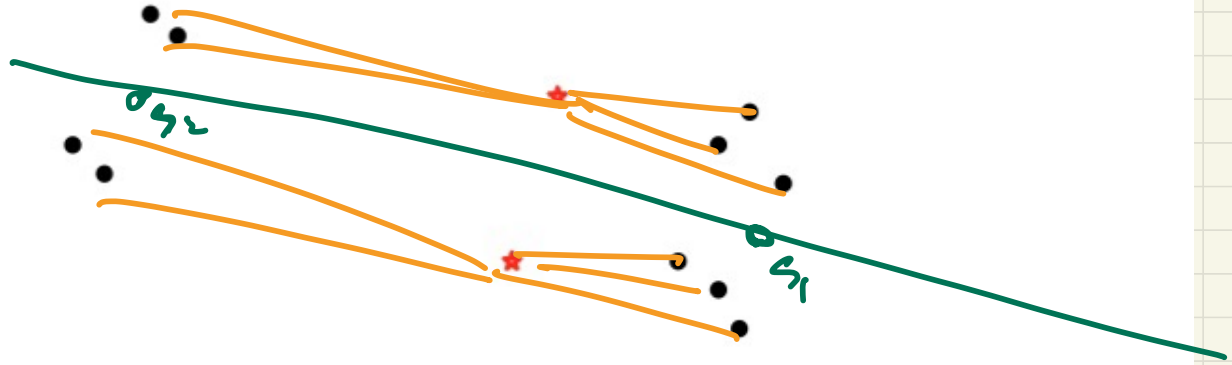


13. Draw the Voronoi diagram of the following set of points.





14. What should you do, if running Lloyd's algorithm for  $k$ -means clustering ( $k = 2$ ), and you reach this scenario, where the algorithm terminates? (The black circles are data points and red stars are the centers).



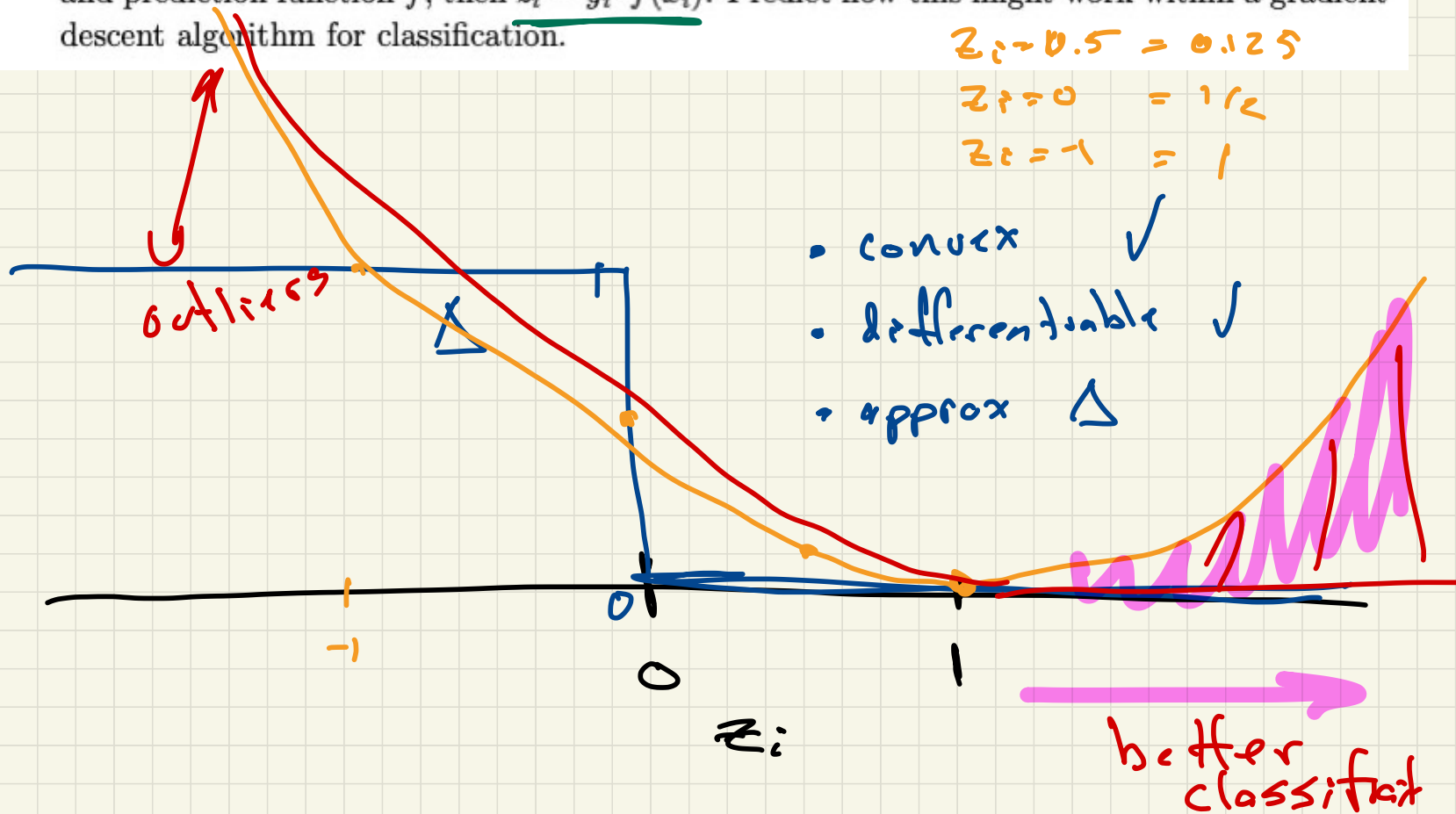
Random Restart

15. Consider the following "loss" function.  $\ell_i(z_i) = (1 - z_i)^2/2$  where for a data point  $(x_i, y_i)$  and prediction function  $f$ , then  $z_i = y_i \cdot f(x_i)$ . Predict how this might work within a gradient descent algorithm for classification.

$$z_i = 0.5 = 0.125$$

$$z_i = 0 = 1/2$$

$$z_i = -1 = 1$$



16. Consider a set of 1-dimensional data points

$$\underline{(x_1 = 0, y_1 = +1)} \quad \underline{(x_2 = 1, y_1 = -1)} \quad \underline{(x_3 = 2, y_1 = +1)} \quad \underline{(x_4 = 4, y_1 = +1)}$$

$$(x_5 = 6, y_1 = -1) \quad (x_6 = 7, y_1 = -1) \quad \underline{(x_7 = 8, y_1 = +1)} \quad (x_8 = 9, y_1 = -1)$$

Predict **-1** or **+1** using a  $k$ NN ( $k$ -nearest neighbor) classifier with  $k = 3$  on the following queries.

- (a)  $x = 3$
- (b)  $x = 9$
- (c)  $x = -1$

