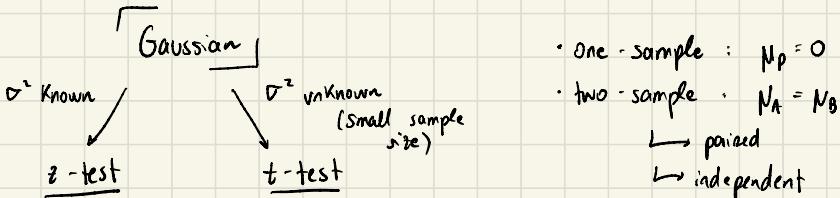



T01 - Statistics

1. a) 1) Choosing the test: → paired samples

- same group of individuals before and after the treatment
(tax increase) \Rightarrow related samples
- Assumption: Quantity of interest is normally distributed.
BUT True variance is unknown.
- small sample size
 - \hookrightarrow paired t-test



$$2) H_0: \bar{N}_d = 0 \quad H_1: \bar{N}_d \neq 0$$

3) Paired t-test:

$$\bar{d} = \frac{1}{10} \cdot \sum_{i=1}^{10} d_i = -8.1 \quad s_d = \sqrt{\frac{1}{n-1} \sum_{i=1}^{10} (x_i - \bar{x})^2}$$

\hookrightarrow Bessel's correction \rightarrow unbiased sample variance

$$T_0 = \frac{\bar{d} - \mu_0}{\sigma_d / \sqrt{n}} = \frac{-8.1}{7.593 / \sqrt{10}} = -3.3734$$

$$4) t^c = t_{1-\alpha/2; n-1}^c = t_{0.975; 9}^c = 2.262$$

p-value: smallest α (significance) at which H_0 can be rejected.

$$5) \text{Rejection criteria} \quad |T_0| > t_{0.975; 9}^c = 2.262$$

$$|T_0| = 3.3734 > 2.262 = t^c \Rightarrow H_0 \text{ can be rejected!}$$

6) Sufficient evidence that a tax increase has an effect on consumption!

2a)

1) Type of test ; independent samples

- independent \rightarrow Welch-test $t_0 = \frac{\bar{x} - \bar{w} - N_0}{\sqrt{s_{\bar{x}-\bar{w}}^2}}$ n t_{df}

$$\Leftrightarrow H_0 = N_D := N_f - N_m \leq 0 \quad ; \quad H_1 = N_D = N_f - N_m > 0$$

3) Test statistics

$$T_0 = \frac{\bar{f} - \bar{m} - N_0}{S_{\bar{f} - \bar{m}}} = \frac{4 - d.5 - 0}{0.769} \approx 1.949$$

$$\bullet \bar{f} = \frac{1}{10} \cdot \sum_{i=1}^{10} f_i = 4$$

$$\left. \begin{aligned} \bullet s_f^2 &= \frac{1}{9} \cdot \sum_{i=1}^{10} (f_i - \bar{m})^2 \approx 3.778 \\ \bullet s_m^2 &\approx 1.714 \end{aligned} \right\} s_{f-m} = \sqrt{\frac{s_f^2}{10} + \frac{s_m^2}{8}}$$

$$\bullet \bar{m} = 2.5$$

$$4) \quad H_0: N_D := N_f - N_m \leq 0 \quad ; \quad H_1: N_D = N_f - N_m > 0$$

$$t^c = t_{1-\alpha}^c; d = t_{0.95; 16}^c = 1.746$$

5) Rejection criteria $|T_0| > t_{0.95, 16}^c = 1.746$

$$T_0 = 1.949 > 1.746 \implies H_0 \text{ can be rejected!}$$

6.) We reject H_0 , which means we acknowledge that

H₂: "On average, women wear their mask longer per day than men."

$$3. \text{ a) } H_0: \mu_{NT} \leq \mu_{OT} \quad H_1: \mu_{NT} > \mu_{OT}$$

we want to prove that NT improves results

$$\Leftrightarrow H_0: N_D := N_{NT} - N_{OT} \leq 0 \quad : \quad H_1: N_D > 0$$

b) $H_0: N_D = N_0 = 0 \Rightarrow$ two-sided (there are two sides (higher and lower) that can violate the hypothesis)

$$H_0: N_D \leq 0 \Rightarrow$$
 One-sided (only one side can violate (greater))

T02 - Regression

1.

a) $\hat{\beta}_1 = \frac{\text{cov}(x, y)}{\text{var}(x)} = \frac{\frac{1}{8}((0.5-3.78)(5-63) + (0.6-3.78)(28-63) + \dots)}{\frac{1}{8}((0.5-3.78)^2 + (0.6-3.78)^2 + \dots)} = \frac{772.300}{133.776} = 5.773 //$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 x = 41.178 //$$

$$\therefore \hat{y}(x) = 41.18 + 5.77x$$

- b) $\hat{\beta}_0$ expected mean value for y when $x = 0$
 $\hat{\beta}_1$ expected change in y per unit of change in x

$\hat{\beta}_0$: A country with a capita gross national product of 0 has (approximately) 41.19% of literate people among the population

$\hat{\beta}_1$: with each increase of \$1.000, the percentage of literate people among the population increases by (approximately) 5.77.

c) Last weeks:

1. Not needed, because only one test is used on regression analysis
2. $H_0: \beta_1 \leq 0$ vs $H_1: \beta_1 > 0$
3. Simply t-test for coefficient. We have

$$SE(\hat{\beta}_1) = \sqrt{\frac{RSS}{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \frac{1}{n-2}} \approx 2.17$$

$$RSS \approx 4,411.4 \quad \sum_{i=1}^n (x_i - \bar{x})^2 \approx 133.77$$

and thus $t_0 \approx 2.66$

4. $\alpha = 0.05$ n - (number of coefficients)
5. For $df = n - 2 = 7$ degrees of freedom, we find a critical value of $t_{9; 0.95}^c = 1.895$ in the t-table.
6. We reject H_0 because $t_0 > t^c$ and conclude that $\hat{\beta}_1$ is statistically significant.

- d) Percentage of literate people (Y) might be over 100 when the capita gross national product (X) gets bigger than 13

We have $X: 13.0 \rightarrow Y: 99$ New Zealand

$$\hat{y} = 41.18 + 5.77x \quad X: 14.0 \rightarrow 121.96$$

T03 - Logistic Regression

$$2. \text{ a) } L(y | x, \beta) = \prod_{i=1}^n p_i^{y_i} (1-p_i)^{1-y_i}$$

$$= p_3 p_4 (1-p_3)(1-p_4)$$

$$5) LL(\beta) = \sum_{i=1}^n y_i \log p_i + (1-y_i) \log (1-p_i)$$

$$LL(\beta) = \log p_3 + \log p_4 + \log (1-p_3) + \log (1-p_4)$$

$$\nabla LL(\beta) = \begin{pmatrix} \frac{\partial LL(\beta)}{\partial \beta_0} \\ \frac{\partial LL(\beta)}{\partial \beta_1} \end{pmatrix}$$

$$\frac{\partial LL}{\partial \beta_i} = \sum_{i=1}^n \frac{\partial LL}{\partial p_i} \cdot \frac{\partial p_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial \beta_j} \quad \text{with} \quad z_i = \beta_0 + \beta_1 x_i \quad \text{and}$$

$$p_i = \sigma(z_i)$$

$$\sigma'(z_i) = \sigma(z_i)(1-\sigma(z_i))$$

$$\frac{\partial LL}{\partial p_i} = \frac{y_i}{p_i} - \frac{1-y_i}{1-p_i} \quad \frac{\partial p_i}{\partial z_i} = \sigma(z_i)(1-\sigma(z_i))$$

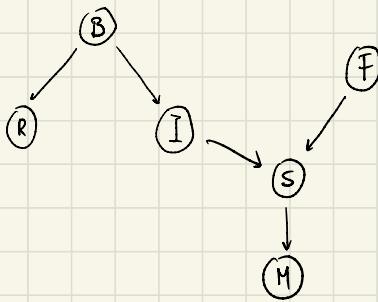
$$\rightarrow \frac{\partial LL}{\partial \beta_0} = \sum_{i=1}^n \left(\frac{y_i}{p_i} - \frac{1-y_i}{1-p_i} \right) \cdot \sigma(z_i) \cdot (1-\sigma(z_i)) \cdot 1$$

$$\frac{\partial LL}{\partial \beta_1} = \sum_{i=1}^n \left(\frac{y_i}{p_i} - \frac{1-y_i}{1-p_i} \right) \cdot \sigma(z_i) \cdot (1-\sigma(z_i)) \cdot x_i$$

c)

T04 - Naive Bayes

1. a)



b) $P(B, F, I, M, R, S) = P(B) \cdot P(F) \cdot P(I|B) \cdot P(R|B) \cdot P(S|I, F) \cdot P(M|S)$

2. Nominal: discr. without order
ordinal: discr. with order

interval: cont. no zero

ratio: cont. with an absolute 0

	cont / disc	range	scale of measurement
age	disc.	(0, +∞)	ratio, ordinal
pref	disc.	{R, B, C}	nominal
money	cont.	(0, +∞)	ratio
product	disc.	{VC, VC, NP}	nominal

T05 - Decision Trees

low since there is a 10/90% division which is good and very pure.

1. a) entropy ((0.1, 0.9)) = $-0.1 \log 0.1 - 0.9 \log 0.9 = 0.469$

b) entropy ((0.8, 0.2)) = $-0.8 \log 0.8 - 0.2 \log 0.2 = 0.722$

c) entropy ((0.5, 0.5)) = $-0.5 \log 0.5 - 0.5 \log 0.5 = 1$ → all classes have the same probability, the node is very impure → high entropy

d) entropy ((0.8, 0.1, 0.1)) = $-0.8 \log 0.8 - 0.1 \log 0.1 - 0.1 \log 0.1 =$

Information gain

$$I(P) = \text{entropy}(p_1, p_2, \dots, p_n) = -p_1 \log_2 p_1 - p_2 \log_2 p_2 - \dots - p_n \log_2 p_n$$

Information

$$I(P) = -(p_1 \log(p_1) + p_2 \log(p_2) + \dots + p_n \log(p_n)) = -\sum_i p_i \log(p_i) = \sum_i p_i \log(1/p_i)$$

Information gain

$$\text{gain}(S, a) = \text{entropy}(S) - \sum_{i \in \text{classes}(a)} \frac{|S_i|}{|S|} \text{entropy}(S_i)$$

Entropy of the parent

Weighted average

$$e) \text{info}([2,3]) = \text{entropy}\left(\frac{2}{5}, \frac{3}{5}\right) = -\frac{2}{5} \log\left(\frac{2}{5}\right) - \frac{3}{5} \log\left(\frac{3}{5}\right) \approx 0.971$$

$$f) \text{info}([5,4]) = \text{entropy}\left(\frac{5}{9}, \frac{4}{9}\right) \approx 0.991$$

$$g) \text{info}([2,3], [5,4]) = \frac{5}{14} \text{info}([2,3]) + \frac{9}{14} \text{info}([5,4]) \approx 0.984$$

\downarrow

0 → pure split

$$h) \text{info}([2,3], [9,0]) = \frac{5}{14} \text{info}([2,3]) + \frac{9}{14} \text{info}([9,0]) \approx 0.347$$

2.

a) Only two possible splits: $x \leq 35$ or $x \leq 37$

35	35	37	40	40	40
F	F	T	F	T	T

$$1. \text{info}([0,2], [3,1]) = 0.541$$

picked!
(purest)

35	35	37	40	40	40
F	F	T	F	T	T

$$2. \text{info}([1,2], [2,1]) \approx 0.918$$

$$\underline{\text{IG}} = \text{info}([3,3]) - 0.541 \geq \text{info}([3,3]) - 0.918$$

\hookrightarrow information of the original set

		1	2	3	
36.8	36.8	37.2	38.3	38.3	39.7
T	F	F	T	F	F

$$\text{info}([2,4]) = 0.918$$

$$1. \text{info}([1,1], [1,3]) \approx 0.874 \quad 2. \text{info}([1,2], [1,2]) \approx 0.918$$

$$3. \text{info}([2,3], [0,1]) \approx 0.809 \quad \text{picked!}$$

∴ Since entropy of the original set is the same, for the highest IG ($\text{IG} = \text{entropy}(S) - \text{entropy}(\text{split})$) we only need the lowest entropy split.

3. a) Past Trend

↑ up ↓ down

$$\text{pos: } p=4 \quad n=2 \quad P: \frac{4}{6} \quad N: \frac{2}{6}$$

$$\text{neg: } p=0 \quad n=0 \quad P: 0 \quad N: 1$$

$$\text{Gini Index: } 1 - P^2 - N^2 = 1 - \frac{4}{9} - \frac{1}{9} = \frac{4}{9}$$

$$\text{Gini Index: } 1 - 0 - 1 = 0$$

$$\therefore \text{Gini Index (past trend)} = \frac{6}{10} \cdot \frac{4}{9} + \frac{4}{10} \cdot 0 = \frac{4}{15}$$

$$\text{Gini Index (open interest)} = \frac{7}{15}$$

$$\text{Gini Index (trading volume)} = \frac{12}{35}$$

b) past trend \rightarrow lower Gini Index \rightarrow purer nodes

T10 - PCA

1. a) $D = \{(-3, -1, -1), (0, -1, 0), (-2, -1, 2), (1, -1, 3)\}$

1. Calculate the means $\bar{d}_1 = \frac{-3+0+(-2)+1}{4} = \frac{-4}{4} = -1$ $\bar{d}_2 = -1$ $\bar{d}_3 = 1$

2. Subtract means

$$X = \begin{bmatrix} -2 & 0 & -2 \\ 1 & 0 & -1 \\ -1 & 0 & 1 \\ 2 & 0 & 2 \end{bmatrix}$$

b) $\text{cov}(x_1, x_1) = \frac{1}{N-1} \sum_{i=1}^N x_{1i}^2 = \frac{1}{3} (4+1+1+4) = \frac{10}{3}$

$$\text{cov}(x_1, x_2) = \frac{1}{3} (0+0+0+0) = 0 = \text{cov}(x_2, x_1)$$

$$\text{cov}(x_1, x_3) = \frac{1}{3} (4-1-1+4) = 2$$

$$\text{cov}(x_2, x_2) = \frac{1}{3} (0) = 0 \quad \text{cov}(x_2, x_3) = 0$$

$$\text{cov}(x_3, x_3) = \frac{1}{3} (4+1+1+4) = \frac{10}{3}$$

$$\Sigma_x = \begin{bmatrix} \text{cov}(x_1, x_1) & \text{cov}(x_1, x_2) & \text{cov}(x_1, x_3) \\ \text{cov}(x_2, x_1) & \text{cov}(x_2, x_2) & \text{cov}(x_2, x_3) \\ \text{cov}(x_3, x_1) & \text{cov}(x_3, x_2) & \text{cov}(x_3, x_3) \end{bmatrix} = \begin{bmatrix} \frac{10}{3} & 0 & 2 \\ 0 & 0 & 0 \\ 2 & 0 & \frac{10}{3} \end{bmatrix}$$

$$c) |\Sigma_x - \lambda I| = 0 \Rightarrow \begin{bmatrix} \frac{10}{3} - \lambda & 0 & 2 \\ 0 & -\lambda & 0 \\ 2 & 0 & \frac{10}{3} - \lambda \end{bmatrix} = 0$$

$$\Leftrightarrow \left(\frac{10}{3} - \lambda\right)(-\lambda)\left(\frac{10}{3} - \lambda\right) + (0 \cdot 0 \cdot 2) + (0 \cdot 0 \cdot 2) - (2 \cdot (-\lambda) \cdot 2) - (0 \cdot 0 \cdot \left(\frac{10}{3} - \lambda\right)) - \left(\left(\frac{10}{3} - \lambda\right) \cdot 0 \cdot 0\right)$$

$$\Leftrightarrow \lambda_1 = 0 \quad \vee \quad \lambda_2 = \frac{4}{3} \quad \vee \quad \lambda_3 = \frac{16}{3}$$

d) Consider all λ :

$$\lambda_1 = 0 : (\Sigma_x - \lambda_1 I) \cdot \vec{v} = 0 \Rightarrow \begin{bmatrix} \frac{10}{3} & 0 & 2 \\ 0 & 0 & 0 \\ 2 & 0 & \frac{10}{3} \end{bmatrix} \cdot \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix} = 0$$

$$\Leftrightarrow \begin{cases} \frac{10}{3} v_1 + 2 v_3 = 0 \\ 2 v_1 + \frac{10}{3} v_3 = 0 \end{cases} \Rightarrow v_1 = v_3 = 0 \quad v_2 \text{ can take any value}$$

$$\text{normalize } |\vec{v}_1| = 1 \Rightarrow \vec{v}_1 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$$

$$\lambda_2 = \frac{4}{3} : \vec{v}_2 = \begin{pmatrix} \frac{\sqrt{2}}{2} \\ 0 \\ -\frac{\sqrt{2}}{2} \end{pmatrix}$$

$$\lambda_3 = \frac{16}{3} : \vec{v}_3 = \begin{pmatrix} \frac{\sqrt{2}}{2} \\ 0 \\ \frac{\sqrt{2}}{2} \end{pmatrix}$$

$$\vec{\Phi} = \begin{pmatrix} \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} & 0 \\ 0 & 0 & 1 \\ \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} & 0 \end{pmatrix}$$

\vec{v}_3 \vec{v}_2

$$\text{For } \lambda_1 = 0 : \frac{0}{0 + \frac{4}{3} + \frac{16}{3}} = 0\%$$

$$\lambda_2 = \frac{4}{3} : \frac{\frac{4}{3}}{0 + \frac{20}{3}} = 20\%$$

$$\lambda = \frac{16}{3} : \frac{\frac{16}{3}}{20} = 80\%$$

$$e) z = x \cdot \vec{\Phi} = \begin{bmatrix} -2 & 0 & -2 \\ 1 & 0 & -1 \\ -1 & 0 & 1 \\ 2 & 0 & 2 \end{bmatrix} \cdot \begin{pmatrix} \frac{\sqrt{2}}{2} \\ 0 \\ \frac{\sqrt{2}}{2} \end{pmatrix} = \begin{bmatrix} -2\sqrt{2} \\ 0 \\ 0 \\ 2\sqrt{2} \end{bmatrix}$$

$$f) z = X \cdot \bar{\Phi} = X \cdot \begin{pmatrix} \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \\ 0 & 0 \\ \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \end{pmatrix} = \begin{pmatrix} -2\sqrt{2} & 0 \\ 0 & \sqrt{2} \\ 0 & -\sqrt{2} \end{pmatrix}$$

2. a) $D \approx z \cdot \bar{\Phi}^T + \text{means}$

$$D \approx \begin{bmatrix} -2\sqrt{2} \\ 0 \\ 0 \\ 2\sqrt{2} \end{bmatrix} \cdot \left(\frac{\sqrt{2}}{2} \ 0 \ \frac{\sqrt{2}}{2} \right) + \begin{bmatrix} -1 & -1 & 1 \\ -1 & -1 & 1 \\ -1 & -1 & 1 \\ -1 & -1 & 1 \end{bmatrix} = \begin{bmatrix} -3 & -1 & -1 \\ -1 & -1 & 1 \\ -1 & -1 & 1 \\ -1 & -1 & 3 \end{bmatrix}$$

$$b) D \approx \begin{bmatrix} -2\sqrt{2} & 0 \\ 0 & \sqrt{2} \\ 0 & -\sqrt{2} \\ 2\sqrt{2} & 0 \end{bmatrix} \cdot \left(\frac{\sqrt{2}}{2} \ 0 \ \frac{\sqrt{2}}{2} \right) + \begin{bmatrix} -1 & -1 & 1 \\ -1 & -1 & 1 \\ -1 & -1 & 1 \\ -1 & -1 & 1 \end{bmatrix} = \begin{bmatrix} -3 & -1 & -1 \\ 0 & -1 & 0 \\ -2 & -1 & 2 \\ 1 & -1 & 3 \end{bmatrix}$$

completely reconstructed since
 $\lambda_3 + \lambda_2$ explain 80% + 20% of the data alone

3. a)
1. Zero Mean Dataset
 2. Covariance
 3. Eigenvalues
 4. Eigenvectors
 5. Projection Matrix

T11 - Convex Optimization

$$1. \quad f(x, y) = a \exp(3x) + \frac{b}{2} xy + y^2$$

$f(x)$ is convex iff $\forall x_1, x_2 \in \text{dom}(f)$

$$\Rightarrow f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2) \quad \forall \lambda \in [0, 1]$$

if f is first-order differentiable, then



$$\Rightarrow f(x_2) \geq f(x_1) + \nabla f(x_1)^T (x_2 - x_1) \quad \forall x_1, x_2 \in \text{dom}(f)$$

f is second-order differentiable $\nabla^2 f(x) \geq 0 \quad \forall x \in \text{dom}(f)$

$$\nabla f(x, y) = \begin{pmatrix} 3a \exp(3x) + \frac{b}{2}y \\ \frac{b}{2}x + 2y \end{pmatrix}$$

$$\nabla^2 f(x, y) = \begin{pmatrix} 9a \exp(3x) & \frac{b}{2} \\ \frac{b}{2} & 2 \end{pmatrix} \geq 0$$

$$\lambda_1, \lambda_2 \geq 0 \quad \lambda_1 + \lambda_2 = 1 \quad \Leftrightarrow \quad \text{Tr}(\nabla^2 f(x, y)) = \lambda_1 + \lambda_2 \geq 0$$

$$\det(\nabla^2 f(x, y)) = \lambda_1 \cdot \lambda_2 \geq 0$$

Note: for every $x, \exp(3x) \geq 0$

$$\text{Tr}(\nabla^2 f(x, y)) = 9a \exp(3x) + 2 \geq 0 \Rightarrow a \geq -\frac{2}{9 \exp(3x)}$$

$$\det(\nabla^2 f(x, y)) = 18 \exp(3x) - \frac{1}{4} b^2 \geq 0 \Rightarrow b \leq \pm \sqrt{4 \times 18 \exp(3x)}$$

$$2. \quad h_2(x) = g_2(y) = g(Ax + b)$$

for $h_2(x)$, $\forall x_1, x_2$ s.t. $Ax_1 + b$ and $Ax_2 + b \in \text{dom}(g_2)$

$$\Rightarrow h_2(\lambda x_1 + (1-\lambda)x_2) = g_2(A(\lambda x_1 + (1-\lambda)x_2) + b)$$

$$= g_2(\lambda(Ax_1 + b) + (1-\lambda)(Ax_2 + b))$$

$$\leq \lambda g_2(Ax_1 + b) + (1-\lambda) g_2(Ax_2 + b)$$

convexity of g_2

$$= \lambda h_1(x_1) + (1-\lambda) h_2(x_2)$$

def of h_2

$$\begin{aligned} \bullet h_1(\lambda x_1 + (1-\lambda)x_2) &= C_1 g_1(\lambda x_1 + (1-\lambda)x_2) + C_2 g_2(\lambda x_1 + (1-\lambda)x_2) \quad \text{convexity of } g_2 \\ &\leq C_1 (\lambda g_1(x_1) + (1-\lambda)g_1(x_2)) + C_2 (\lambda g_2(x_1) + (1-\lambda)g_2(x_2)) \\ &= \lambda [C_1 g_1(x_1) + C_2 g_2(x_1)] + (1-\lambda)[C_1 g_1(x_2) + C_2 g_2(x_2)] \\ &= \lambda h_1(x_1) + (1-\lambda)h_1(x_2) \end{aligned}$$

$$\begin{aligned} \bullet h_3(\lambda x_1 + (1-\lambda)x_2) &= \max \{ g_1(\lambda x_1 + (1-\lambda)x_2), g_2(\lambda x_1 + (1-\lambda)x_2) \} \\ &\leq \max \{ \lambda g_1(x_1) + (1-\lambda)g_1(x_2), \lambda g_2(x_1) + (1-\lambda)g_2(x_2) \} \\ &\leq \max \{ \lambda h_3(x_1) + (1-\lambda)h_3(x_2), \lambda h_3(x_2) + (1-\lambda)h_3(x_1) \} \end{aligned}$$

$g_1(x) \leq h_3(x) \quad \swarrow$
 $g_2(x) \leq h_3(x) \quad \searrow$

$$3. f(x,y) = 2x^2 + 0.5y^2 - 3x - y - 2xy + 5 \quad z^{(1)} = (0,0)$$

$$\nabla f = \begin{pmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{pmatrix} = \begin{pmatrix} 4x - 2y - 3 \\ y - 2x - 1 \end{pmatrix}$$

$$\boxed{\text{line-search: } \alpha^{(n)} : \underset{\alpha}{\operatorname{argmin}} f(x^{(n)} - \alpha \cdot \nabla f(x^{(n)}))}$$

↳ minimizes the result of f in order to iteratively approach a local minimum

$$\nabla f(0,0) = \begin{pmatrix} -3 \\ -1 \end{pmatrix} \quad f(0,0) = 5$$

$$\begin{aligned} \alpha^{(1)} &= \underset{\alpha}{\operatorname{argmin}} f(z^{(1)} - \alpha \cdot \nabla f(z^{(1)})) \\ &= \underset{\alpha}{\operatorname{argmin}} 12.5\alpha^2 - 10\alpha + 5 \\ &= \frac{10}{25} = 0.4 \end{aligned}$$

$$② z^{(2)} = z^{(1)} - \alpha^{(1)} \cdot \nabla f(z^{(1)}) = \begin{pmatrix} 1.2 \\ 0.4 \end{pmatrix}$$

$$\begin{aligned} \alpha^{(2)} &= \underset{\alpha}{\operatorname{argmin}} f(z^{(2)} - \alpha \cdot \nabla f(z^{(2)})) \\ &= \underset{\alpha}{\operatorname{argmin}} 12.5\alpha^2 - 10\alpha + 3 \\ &= \frac{10}{25} = 0.4 \end{aligned}$$

$$③ z^{(3)} = z^{(2)} - \alpha^{(2)} \cdot \nabla f(z^{(2)}) = \begin{pmatrix} 1.2 \\ 0.4 \end{pmatrix} - 0.4 \cdot \begin{pmatrix} 1 \\ -3 \end{pmatrix} = \begin{pmatrix} 0.8 \\ 1.6 \end{pmatrix}$$

T12 - Neural Networks

1. 1. $\cdot K=1 \Rightarrow F(x) = W_1 x \rightarrow$ linear with no bias TRUE

$$\begin{aligned} \cdot K=2 \quad \nabla l(F) &= \frac{\partial l}{\partial F} \cdot \frac{\partial F}{\partial x} = -\frac{2}{n} \sum_{i=1}^n (y_i - F(x_i)) \cdot (-W_1 W_2) \\ &= (W_2 W_1) \cdot \frac{2}{n} \sum_{i=1}^n (y_i - F(x_i)) \end{aligned}$$

TRUE

- Increasing $K \Rightarrow F(x) = W_K W_{K-1} \dots W_1 x$
 - \hookrightarrow all linear transformations can be mapped to $W_J = W_K W_{K-1} \dots W_1$

- $W_K \in \mathbb{R}^{d_1 \times d_2}$ ($d_1, d_2 > 1$) leads to output of size $d_2 \times K$ for a K that depends on the previous results size

$F(x) \in \mathbb{R}^{1 \times 1}$ since it outputs a scalar label
 $\therefore d_1 > 1$ but needs to be 1 FALSE

2. $K=3$ W_3 scalar

$$F(x) = W_3 W_2 W_1 x \quad l_x = (y - F(x))^2$$

$$\frac{\partial l_x}{\partial W_3} = \frac{\partial l}{\partial F} \cdot \frac{\partial F}{\partial W_3} = -2(y - F(x)) \cdot W_2 W_1 x$$

$$2. \text{ a)} \quad x \in \mathbb{R}^{2 \times 1} \quad W^{[1]} \in \mathbb{R}^{2 \times 2} \quad W^{[2]} \in \mathbb{R}^{2 \times 2} \\ b^{[2]} \in \mathbb{R}^{1 \times 2} \quad b^{[1]} \in \mathbb{R}^{1 \times 2}$$

Trainable parameters: $2 \times 2 + 2 + 2 \times 1 + 1 = 9$

$$b) \hat{y} = \sigma(X^{[1]} \cdot W^{[1]} + \vec{1} \times b^{[1]}) = \sigma(\sigma(X \cdot W^{[2]} + \vec{1} \times b^{[2]})) W^{[1]} + \vec{1} \times b^{[1]}$$

$b^{[2]} \in \mathbb{R}^{1 \times 2} \longrightarrow \vec{1} \times b^{[2]} \in \mathbb{R}^{n \times 2} \quad [\dots] \rightarrow \begin{bmatrix} \dots \\ \dots \end{bmatrix}$

$$c) \quad z^{(1)} = X \cdot W^{(1)} + b^{(1)}$$

$$= \begin{pmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{pmatrix}$$

$$X^{(1)} = \sigma(z^{(1)}) = \begin{pmatrix} 0.5 & 0.5 \\ 0.5 & 0.931 \\ 0.931 & 0.5 \\ 0.931 & 0.931 \end{pmatrix}$$

$$\sigma = \frac{1}{1+e^{-x}}$$

$$z^{(2)} = X^{(1)} \cdot W^{(2)} + b^{(2)} = \begin{pmatrix} 0.5 & 0.5 \\ 0.5 & 0.931 \\ 0.931 & 0.5 \\ 0.931 & 0.931 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ -1 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ -0.231 \\ 0.231 \\ 0 \end{pmatrix}$$

$$X^{(2)} = \sigma(z^{(2)}) = \begin{pmatrix} 0.5 \\ 0.443 \\ 0.557 \\ 0.5 \end{pmatrix} = \hat{y}$$

$$y = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

$$\therefore L(y, \hat{y}) = \frac{1}{4} \sum_{i=1}^4 l(Y_i, \hat{Y}_i) = 0.696$$

$$d) \quad \frac{\partial L}{\partial w^{(2)}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial w^{(2)}}$$

$$\frac{\partial L}{\partial w^{(1)}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z^{(1)}} \cdot \frac{\partial z^{(1)}}{\partial x^{(1)}} \cdot \frac{\partial x^{(1)}}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial w^{(1)}}$$

$$\frac{\partial L}{\partial b^{(2)}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial b^{(2)}}$$

$$\frac{\partial L}{\partial b^{(1)}} = " \cdot " \cdot " \cdot " \cdot " \cdot \frac{\partial z^{(1)}}{\partial b^{(2)}}$$

