Question 1

Linear Regression:

1)
$Y$ ⟶ output     $(n \times 1)$
$X$ ⟶ input      $(n \times P)$
$\beta$ ⟶ coefficients   $(p \times 1)$
$\varepsilon$ ⟶ errors     $(n \times 1)$

$$Y = X\beta + \varepsilon$$

Assuming that the output is a linear model of the input & errors have mean zero, & constant variance

2) OLS — $\overset{sum \theta}{\underset{of}{}}$ squares of (actual − predicted)

$$\sum (y_i - \hat{y}_i)^2 \Longleftarrow$$

MSE — OLS / number of outputs

$$MSE = \frac{1}{n} \left( \| y - X\beta \|^2 \right)$$

actual

$$= \frac{1}{n} (y - X\beta)^T (y - X\beta)$$

3) $\dfrac{\sum (OLS)}{\partial \beta}$

3) $\dfrac{\partial}{\partial \beta} (OLS / MSE) =$

$$OLS = (y^T - \beta^T x^T)(y - X\beta)$$

Key step — recognizing the dimension of $y^T x \beta \longrightarrow (1 \times 1)$ $\therefore (y^T x \beta)^T = y^T x \beta$

$$y^T y - y^T x \beta - \beta^T x^T y + \beta^T x^T x \beta$$

$$y^T y - 2 \beta^T x^T y + \beta^T x^T x \beta$$

$$\frac{d}{d\beta} \left( y^T y - 2 \beta^T x^T y + \beta^T x^T x \beta \right)$$

$$0 = 0 - 2 x^T y + \otimes x^T x \hat{\beta} + (x^T x)^T \hat{\beta}$$

$$\downarrow$$

$$\frac{d (\beta^T \alpha)}{d\beta} = \alpha$$

$$2 x^T y = 2 x^T x \hat{\beta}$$
$$\hat{\beta} = (x^T x)^{-1} x^T y$$

4) $\det(x^T x) \neq 0$ & $x^T x$ should be a square matrix. $x^T x$ will always be square. If we have multicollinearity occurs then we could have then transformed the matrix to make a column zero and thus |det| would be zero

5) $\hat{y} = x\hat{\beta} + \varepsilon$

$$y = x\beta + \varepsilon$$
$$\boxed{y - \hat{y} = \varepsilon}$$

To prove : $x^T (y - \hat{y}) = 0$

$$= x^T (y - x\hat{\beta})$$

$$= x^T y - x^T x \hat{\beta}$$
$$= 0 \quad (\text{Above derivation of } \hat{\beta})$$

6) $\quad J(\beta) = \dfrac{1}{2n} \| x\beta - y \|^2$

$$\nabla_\beta (x\beta - y)^T (x\beta - y)$$

$\dfrac{\partial}{\partial \beta}$

$$\dfrac{\partial}{\partial \beta} (\beta^T x^T - y^T)(x\beta - y)$$

$$\cancel{\beta^T x^T x} \cancel{y^T x\beta}$$

$$\beta^T x^T x \beta - y^T x\beta = \beta^T x^T y + y^T y$$

$$\boxed{\dfrac{\partial}{\partial \beta} = \dfrac{-2x^T y + 2x^T x\beta}{2n}}$$

$$\dfrac{\partial}{\partial \beta} = \dfrac{x^T x\beta - x^T y}{n}$$

By substituting this value we will get the result.

Batch gradient rule

$$\beta^{(t+1)} = \beta^{(t)} - \alpha \nabla_\beta J(\beta^{(t)})$$

10) Heavy tails — extreme value exist in a large number thus causing difficulty to predict the extreme outliers & error will be misleading

∞. V shaped — upward — right-skewed & positive outliers

downward — left skewed & negative outliers

12) $E\left((y - \hat{f}(x))^2\right) = Bias^2 + Var + \sigma^2$

$y = f(x) + \varepsilon$
$E(\varepsilon) = 0$
$Var(\varepsilon) = E(\varepsilon^2) = \sigma^2$

$E\left(\left((f(x) - \hat{f}(x)) + \varepsilon\right)^2\right)$

$= E\left((f(x) - \hat{f}(x))^2 + 2(f(x) - \hat{f}(x))\varepsilon + \varepsilon^2\right)$

$E\left(f^2(x) - 2f(x)\hat{f}(x) + \hat{f}(x)\right) + \sigma^2 + \sigma^2$

$E\left(\left((f(x) - E[\hat{f}(x)]) - (\hat{f}(x) - E\hat{f}(x))\right)^2\right)$

$E\left((f(x) - E(\hat{f}(x)))^2 - 2(f(x) - E(\hat{f}(x)))(\hat{f}(x) - E(\hat{f}(x)))\right.$

$\left. + (\hat{f}(x) - E(\hat{f}(x)))^2\right) + \sigma^2$

$E\left(f(x) - E(\hat{f}(x))\right)^2 \qquad E\left((\hat{f}(x) - E(\hat{f}(x)))^2\right)$
$+ bias^2 \qquad\qquad + \qquad\qquad var$

$Bias^2 + Var + \sigma^2$

(3)  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$

$y = \alpha_0 + \alpha_1 x_1 + u$

a)  $E(\alpha_1) = \dfrac{Cov(x_1, y)}{Var(x_1)}$

$= \dfrac{Cov(x_1, \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon)}{Var(x_1)}$

$= Cov(x_1, \beta_0) + Cov(x_1, \beta_1 x_1) + Cov(x_1, \beta_2 x_2) + Cov(x_1, \varepsilon)$

$= \dfrac{\beta_1 Var(x_1) + \beta_2 Cov(x_1, x_2)}{Var(x_1)}$

$= \beta_1 + \beta_2 \dfrac{Cov(x_1, x_2)}{Var(x_1)}$

b)  $\left( f(\alpha_1) - \beta_1 \right) x_2$ shows us the bias of $\alpha_1$ towards

c)  for removing the the bias either $\beta_2$ or $Cov(x_1, x_2)$ should be zero.

14) ~~a) We can check for the ratio of coefficients of each column if the valu~~

a) We can check for the ratio of max & min eigen value

if the ratio ~~is~~ close to one~~m~~ the no multicollinearity otherwise if values are large the multicollinearity occurs.

b)

$$Var(\hat{\beta_j}) = \frac{\sigma^2}{(n-1)\, Var(n_j)} \cdot \boxed{\frac{1}{1-R_j^2}} \longrightarrow \text{Variance inflation factor}$$

correlation

so $Var(\hat{\beta_j}) \uparrow$ as ~~corretion~~ correlation $\uparrow$

c) Instability is because ~~the~~ error of the loss functions are flattened

As long as the matrix in invertible bias will be singular

Question 2

2) for education level to handle the missing data ~~~
   (mode) will be the most suitable as mean & median
   of a category cannot be calculated . Same thing
   for city
   Salaries median would be suitable as this property
   is non-linear and thus it would be able to
   manage the outliers well. The same logic could
   work for performance score

3) Encode categorical features

   assigning educational levels numbers as there is a clear
   hierarchy in the levels where PhD is the highest and
   high school is lowest

   for remote worker either 0 or 1

   one hot encoding for the rest of the features
   because assigning numbers will be a problem
   as by assigning numbers @ it will generate an
   ordinality between the categories.

4) for stratification the variable used is salary.
   Salary helps us bifurcate the data so that extreme
   high outliers does not affect the low salaries as
   would have happened if the datasets were randomly
   chosen.