1) $SSE(\vec{W}) =$ Cost function, the squared sum error

$E(\vec{W}) = \sum\limits_{i=1}^{N} (y^{(i)} - \hat{y}^{(i)})^2$

where, $\hat{y} = f(\vec{x}^{(i)}, \vec{W})$

$\Rightarrow E(\vec{W}) = (\vec{y} - x\vec{W})^T (\vec{y} - x\vec{W})$

$\vec{W}^* = \underset{\vec{W}}{\arg\min}(E(\vec{W}))$

$\Rightarrow \nabla E(\vec{W}) = 0 \quad\quad —①$

$E(\vec{W}) = \vec{y}^T\vec{y} - \vec{W}^T x^T \vec{y} - \vec{y}^T x\vec{W} + \vec{W}^T x^T x\vec{W}$

$\nabla E(\vec{W}) = \cancel{0\cancel{x\cancel{y}}} 0 - 2x^T\vec{y} + 2x^T x \vec{W}$

From ①,

$-2x^T\vec{y} + 2x^T x \vec{W}$

$\Rightarrow \boxed{\vec{W}^* = (x^T x)^{-1} x^T y}$

$\boxed{\hat{y} = x\vec{W}^*}$

Also, $\nabla^2 E(\vec{W})$ needs to be positive definite

$\cancel{@} \nabla^2 E(\vec{W}) = 2x^T x$

2). Using Basis functions $\phi(x)$.

$$X = \begin{bmatrix} x_0^1 & x_1^1 & \cdots & x_d^1 \\ x_0^2 & & & \\ \vdots & & & \\ x_0^N & - - - & - & x_d^N \end{bmatrix}_{N \times (d+1)} \quad \vec{Y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}_{N \times 1}.$$

Let the basis function be

$$\phi(X) = \begin{bmatrix} \phi_0^1[x] & \phi_1^1(x) \cdots & \phi_{d}^1(x) \\ \phi_0^2(x) & & \\ \vdots & & \\ \phi_0^N(x) & - - - - & \phi_{d}^N(x) \end{bmatrix}$$

Here, we $\phi_0^{(i)}(x) = 1$ for bias.

Squared Sum Error $\rightarrow SSE(\vec{w})$

$$SSE(\vec{w}) = |\vec{Y} - \hat{y}|^2 = |\vec{Y} - \phi(x)\vec{w}|^2$$

Differentiating and equating to zero, we get

$$\nabla SSE(\vec{w}) = 0$$

$$\Rightarrow \nabla |\vec{Y} - \phi(x)\vec{w}|^2 = 0$$

$$\Rightarrow \boxed{\vec{w} = (\phi^T(x)\,\phi(x))^{-1}\,\phi^T(x)\vec{Y}}$$

This can be solved similarly to the first problem.

3) Given,

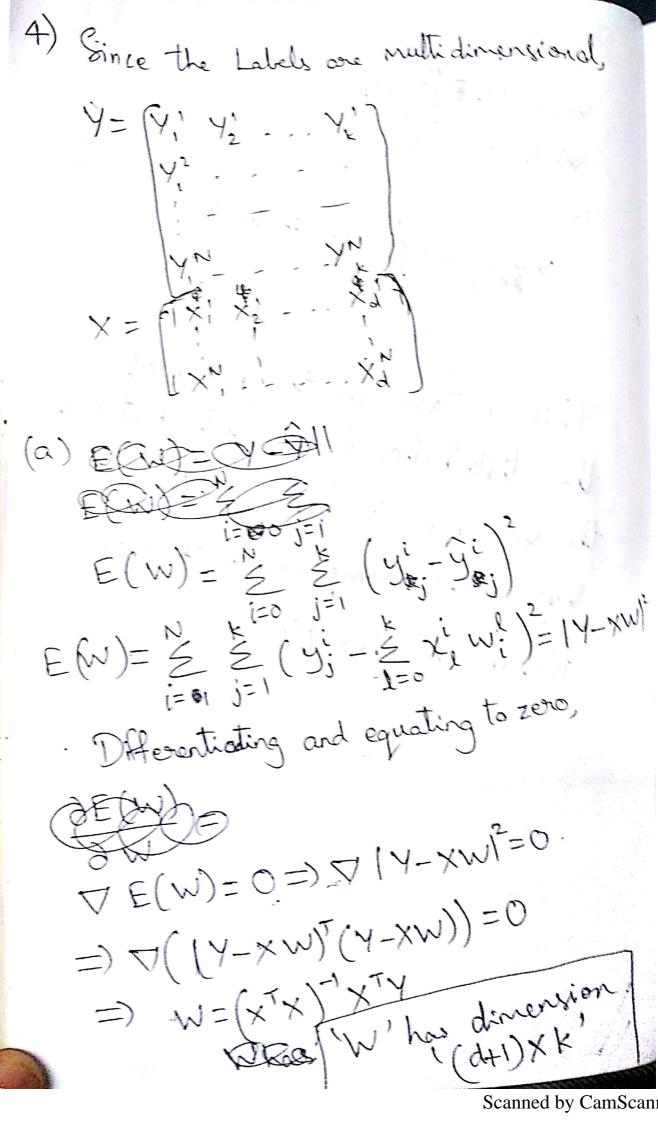$$\sigma(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{e^x + 1}$$

We know that,

$$Tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{e^{2x} - 1}{e^{2x} + 1} = \frac{e^{2x}}{e^{2x} + 1} - \frac{1}{e^{2x} + 1}$$

$$= \frac{e^{2x}}{e^{2x} + 1} - \left(\frac{1 + e^{2x} - e^{2x}}{e^{2x} + 1}\right) = 2\sigma(2x) - 1$$

$$\boxed{\therefore \quad Tanh(x) = 2\sigma(2x) - 1}$$

$$\hat{y}(x, \vec{\omega}) = W_0 + \sum_{j=1}^{M} w_j \sigma\left(\frac{x - \mu_j}{S}\right)$$

$$\hat{y}(x, \vec{\mu}) = U_0 + \sum_{j=1}^{M} u_j \left(2\sigma\left(2\left(\frac{x - \mu_j}{S}\right)\right) - 1\right)$$

$$= \left[\mu_0 - \sum_{j=1}^{M} u_j\right] \oplus + \sum_{j=1}^{M} u_j 2\sigma\left(\frac{2x - 2\mu_j}{S}\right)$$

By comparision,

$$W_0 = U_0 \neq \sum_{j=1}^{M} U_j$$

⊗ $W_1 = U_1, W_2 = U_2 \cdots \cdots$

(or)

$$W_k = \begin{cases} U_0 - \sum_{j=1}^{M} U_j & \text{if } k = 0 \\ \\ U_k & \text{else} \end{cases}$$

This is true as the mean of the distribution is constant but only the varience is changing. Hence, $W_{0j}$ can be compared with $U_j$.

Problem number ⑦.

4) Since the Labels are multidimensional,

$$\dot{Y} = \begin{pmatrix} Y_1^1 & Y_2^1 & \cdots & Y_k^1 \\ Y_1^2 & & & \\ \vdots & & & \\ Y_1^N & \cdots & & Y_k^N \end{pmatrix}$$

$$X = \begin{pmatrix} 1 & X_1^1 & X_2^1 & \cdots & X_d^1 \\ \vdots & & & & \\ 1 & X_1^N & \cdots & & X_d^N \end{pmatrix}$$

(a) $E(w) = \| Y - \hat{Y} \|$

$$E(w) = \sum_{i=0}^{N} \sum_{j=1}^{k} \left( Y_{ij}^i - \hat{Y}_{ij}^i \right)^2$$

$$E(w) = \sum_{i=1}^{N} \sum_{j=1}^{k} \left( y_j^i - \sum_{l=0}^{k} x_l^i w_l^i \right)^2 = |Y - Xw|^2$$

. Differentiating and equating to zero,

$$\frac{\partial E(w)}{\partial w} = 0$$

$$\nabla E(w) = 0 \Rightarrow \nabla |Y - Xw|^2 = 0.$$

$$\Rightarrow \nabla \left( (Y - Xw)^T (Y - Xw) \right) = 0$$

$$\Rightarrow w = (X^T X)^{-1} X^T Y$$

where $w$ has dimension $(d+1) \times k$

(b) Similarly,

$$\phi(X) = \begin{bmatrix} \phi_0^1(x) & \phi_1^1(x) & \ldots & ; & \phi_d^1(x) \\ \vdots & & & & \\ \phi^N(x) & - & - & \cdots & \phi_d^N(x) \end{bmatrix}_{N \times (d+1)}$$

$$\overline{W} = (\phi^T \phi)^{-1} \cdot \phi^T y \quad \text{as the parameters}$$

'$X$' are replaced with their basis functions
'$\phi(x)$'.

5) Given,

$$E(\vec{w}) = \sum_{i=1}^{N} r_i (y^{(i)} - \hat{g}^{(i)})^2 = \sum_{i=1}^{N} r_i \left( y^{(i)} - \sum_{j=0}^{d} x_j^{(i)} w_j \right)^2$$

$$E(\vec{w}) \equiv \sum_{i=0}^{N} \left( \left( \sqrt{r_i}\, y^{(i)} \right) - \sum_{j=0}^{d} \left( x_j^{(i)} \sqrt{r_i} \right) w_j \right)^2.$$

This form is similar to the Sum of squared
Error (SSE) form where the labels and inputs
are

$$\vec{y}' = \begin{bmatrix} y^1 \sqrt{r_1} \\ y^2 \sqrt{r_2} \\ \vdots \\ y^N \sqrt{r_N} \end{bmatrix} = \begin{bmatrix} \sqrt{r_1} & 0 & 0 & \cdots & 0 \\ 0 & \sqrt{r_2} & \cdots & & 0 \\ \vdots & & - & - & \sqrt{r_N} \\ 0 & - & - & - & \end{bmatrix} \begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^N \end{bmatrix}$$

$$x' = \begin{bmatrix} \sqrt{r_1} & 0 & \cdots & & 0 \\ 0 & \sqrt{r_2} & - & - & 0 \\ 0 & - & - & - & \sqrt{r_N} \end{bmatrix} \begin{bmatrix} x_0^1 & x_1^1 & \cdots & x_d^1 \\ \vdots & & & \vdots \\ x_0^N & - & - & x_d^N \end{bmatrix}$$

Let $R = \begin{bmatrix} \sqrt{r_1} & 0 & \cdots & & 0 \\ 0 & \sqrt{r_2} & \cdots & & 0 \\ \vdots & & & & \\ 0 & \cdots & & & \sqrt{r_N} \end{bmatrix}$

$$\therefore E(\vec{w}) = (\vec{Y'} - x'w)^2$$

where, $x' = RX$, $\vec{Y'} = R\vec{Y}$

Minimising the cost function,

$$\nabla E(\vec{w'}) = 0$$

$$\Rightarrow \nabla (\vec{Y'} - x'\vec{w})^2 = 0$$

$$\Rightarrow \vec{w} = ((x')^T x')^{-1} (x')^T \vec{y'}$$

$$= ((RX)^T (RX))^{-1} (RX)^T (RY)$$

$$\boxed{\vec{w} = (x^T R^T R X)^{-1} (x^T R^T R Y)}$$

This is the optimum set of weights for

a weighted sum squared error case.

6) SSE cost function with L2 regularisation is

$$E(\vec{w}) = |\vec{Y} - X\vec{W}|^2 + \lambda W^T W$$

$\lambda \rightarrow$ Lagrangian parameter

**Sol** Minimising the cost, we get

$$\nabla E(\vec{W}) = \nabla (|\vec{Y} - X\vec{W}|^2 + \lambda W^T W) = 0$$

$$\Rightarrow -2X^T(\vec{Y} - X\vec{W}) + 2\lambda I \vec{W} = 0$$

$$\Rightarrow \vec{W} = (X^T X + \lambda I)^{-1} X^T \vec{Y}, \text{ where}$$

X is the set of features not including

**Uses of regularisation:**

- **Overfitting** - It means the model has more parameters than can be justified by the data. This results in the models capturing the irregularities and the noise in the data. Simply put, an overfitted model is too good to be true.

→ One way to overcome overfitting is regularisation. This significantly reduces the variance of the model without affecting the ~~higher bias~~ substantial increase in the bias.

→ By increasing the $\lambda$, it prevents overfitting but increasing it beyond a certain value, results in ~~improper~~ higher error rates.

1) Given,

$$x' = x + \text{Noise}$$

$$E(N_{xy}) = 0, \quad E(N_{xy}^2) = \sigma^2$$

$$SSE(\vec{w}) = \sum_{i=1}^{N} \left[ y^i - \sum_{j=0}^{d} x_j^i w_j \right]^2, \text{ here.}$$

Adding the noise, we get :

$$SSE'(\vec{w}) = \sum_{i=1}^{N} \left[ y^i - \sum_{j=0}^{d} (x_j^i + N_j^i) w_j \right]^2$$

$$= \sum_{i=1}^{N} \left[ \underbrace{y^i - \sum_{j=0}^{d} x_j^i w_j}_{a} - \underbrace{\sum_{j=0}^{d} N_j^i w_j}_{b} \right]^2$$

$$E[SSE'(\vec{w})] = \sum_{i=1}^{N} E[a^2] + E[b^2] - E[2ab]$$

$$\boxed{NE[a^2] = E[SSE(\vec{w})]}$$

$$E[b^2] = E\left[ \left( \sum_{j=0}^{d} N_j^i w_j \right)^2 \right]$$

$$= E\left[ (N_{0}^i w_1)^2 + (N_2^i w_2)^2 + \dots (N_d^i w_d)^2 \right.$$
$$\left. + (N_1^i w_1 N_2^i w_2) + \dots \right]$$

$$\longrightarrow 0$$

$$\boxed{E[N_i N_j] = \delta_{ij}\sigma^2}$$

$$= E\left[ (N_1^i w_1)^2 + \dots (N_d^i w_d)^2 \right]$$

$$= \left[ E[(N_1^i w_1)^2] + E[(N_2^i w_2)^2] \dots + E((N_d^i w_d)^2) \right]$$

~~E[SSE(W)]~~ $\boxed{\sigma^2 \left[\sum_{j=01}^{d} W_j^2\right] = E[b^2]}$

~~E[2ab] = E[2b(y)]~~

$E(2ab) = E\left[2 Err(\vec{W})\left(\sum_{j=01}^{d} N_j^i W_j\right)\right]$

Here, $Err(\vec{W})$ is constant, $W_j$ are consta~~nt~~

$E[2ab] = 2 Err(\vec{W}) E\left[N_1^i W_1 + N_2^i W_2 \cdots + N_d^i W_d\right.$

$= 2 Err(\vec{W})\left[E(N_1^i)W_1 + E(N_2^i)W_2 \cdots\right.$

$= 0$  $\boxed{\therefore E[Noise] = 0}$

$\therefore E[SSE'(\vec{W})] = E[SSE(\vec{W})] + \sigma^2\left[\sum_{j=0}^{d} W_j^2\right]$

$= |Y - X\vec{W}|^2 + \sigma^2 W^T W$

This is simi~~la~~lar to L2 norm where

$\boxed{\lambda = \sigma^2}$

8) Maximum aposterior expression

$$P(W/X,Y,\alpha,\beta) \propto P(Y/X,W,\alpha,\beta)P(W/\alpha)$$

$\hookrightarrow$ from Bayes theorem

~~We know~~ Given,

$$\vec{Y} \sim N(X\vec{W}, \sigma^2 I)$$

$$\Rightarrow P(Y/X,W,\alpha,\beta) \sim \prod_{i=1}^{N} N(\hat{Y}, \sigma^2 I)$$

$$P(W/\alpha) \sim N(0, \alpha^2 I)$$

$$P(W/X,Y,\alpha,\beta) \propto \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^N \exp\left(\frac{(Y-\hat{Y})^T(Y-\hat{Y})}{2\sigma^{-2}}\right)\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^d$$

$$\exp\left(\frac{-W^TW}{2\alpha^2}\right)$$

$$\propto \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^{N+d} \exp\left(-\left(\frac{(Y-\hat{Y})^T(Y-\hat{Y})}{2\sigma^2} + \frac{W^TW}{2\alpha^2}\right)\right)$$

$$\Rightarrow P(W/X, Y, \alpha, \beta) = C \exp\left(-\left[\frac{(y-\hat{y})^T(y-\hat{y})}{2\sigma^2} + \frac{W^TW}{2\sigma^2}\right]\right)$$

for Maximising $P(W/X, Y, \alpha, \beta)$, since the exponent is negative, because maximising the power increases the exponent.

$$\underset{W}{\text{ArgMin}}\left(\frac{(y-\hat{y})^T(y-\hat{y})}{2\sigma^2} + \frac{W^TW}{2\alpha^2}\right) = \underset{W}{\text{Argmax}} P(W_{|})$$

$$\Rightarrow \nabla\left(\frac{|y-\hat{y}|^2}{2\sigma^2} + \frac{|W|^2}{2\alpha^2}\right) = 0$$

$$\Rightarrow \nabla\left(|y-\hat{y}|^2 + \frac{\sigma^2}{\alpha^2}|W|^2\right) = 0$$

This is similar to the ridge regression

where $\lambda = \frac{\sigma^2}{\alpha^2}$