$$Y = XW + E$$
$$N \times D \quad D \times 1$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

$\leftarrow$ features $\rightarrow$

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix} \Big\uparrow \text{sample}$$
$$N \times D$$

$N = $ all the sample

$$E = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_N \end{bmatrix} \leftarrow \text{error} \\ \text{residue}$$

**side**

$f(x) = b^T x$

$f'(x) = b$

b/c

if $f(x) \in \mathbb{R}$

$x \in \mathbb{R}^D$
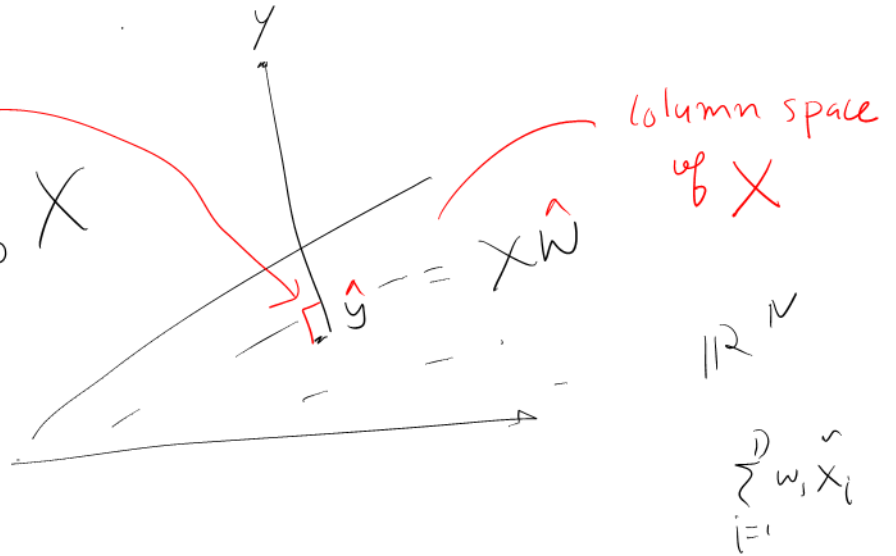
$$f'(x) = \begin{bmatrix} \dfrac{\partial f(x)}{\partial x_1} \\ \vdots \\ \dfrac{\partial f(x)}{\partial x_D} \end{bmatrix}$$

$$RSS = \sum_{i=1}^{N} (y_i - x_i^T w)^2$$
$$= \sum_{i=1}^{N} \epsilon_i^2 = \| E \|_2^2$$

want to find

$\hat{w}$ s.t $\hat{y_i} = x_i^T \hat{w}$

minimizes $\sum_{i=1}^{N} (y_i - \hat{y_i})^2$

$\Downarrow$

minimize $(Y - XW)^T (Y - XW)$

or in Matrix notation

$(Y - X\hat{w})$

$\hat{y} = X\hat{w}$ s.t $(y - \hat{y})^T (y - \hat{y})$ is minimum

$$= \bar{y} = XW$$
$$N \times D \quad D \times 1$$

Let search for such $\hat{w}$. Note the action of a generic $W \in \mathbb{R}^D$ on $X$

$$= \begin{bmatrix} \tilde{x}_1 & \tilde{x}_2 & \cdots & \tilde{x}_D \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix}$$

$\tilde{x}_j \in \mathbb{R}^D$ is a column of $X$

ie $\bar{y} = XW = w_1 \tilde{x}_1 + w_2 \tilde{x}_2 + \cdots w_D \tilde{x}_D \in \mathbb{R}^N$

weighted linear combination of the columns of the matrix $X$, where weights $w_i$ comes from vector $W \in \mathbb{R}^D$

Hence as $w$ varies it spans column space of matrix $X$

Our desired vector $\hat{y} = X\hat{w}$ will also be a point in this column space of $X$ which minimize $(Y-XW)^T(Y-XW)$

ie $w = \hat{w}$

If $\hat{y} = X\hat{w}$ is such a point then

vector $Y-\hat{y}$ has to be perpendicular to column space of $X$

OR

$Y - \hat{y}$ has to be perpendicular to columns of $X$

OR



y

column space of $X$

$\mathbb{R}^N$

$\sum_{i=1}^{D} w_i \hat{x}_i$

side $a \perp b \in \mathbb{R}^N$

$\Longleftrightarrow$

$a^T b = b^T a$

$= \langle a, b \rangle = 0$

This a algebraic definition of two vector $a, b$ being perpendicular is

$\tilde{x}_j^T (y - \hat{u}) = 0 \quad \forall j = 1 : D$

(all)

$x_j^T (y - X\hat{w}) = 0$

$X^T (y - X\hat{w}) = 0$

$X^T y - X^T X \hat{w} = 0$

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

Ridge Regression

MLE can overfit. It tries to explain current evidence, Not good for noisy situations.

$$P(D \mid W)$$

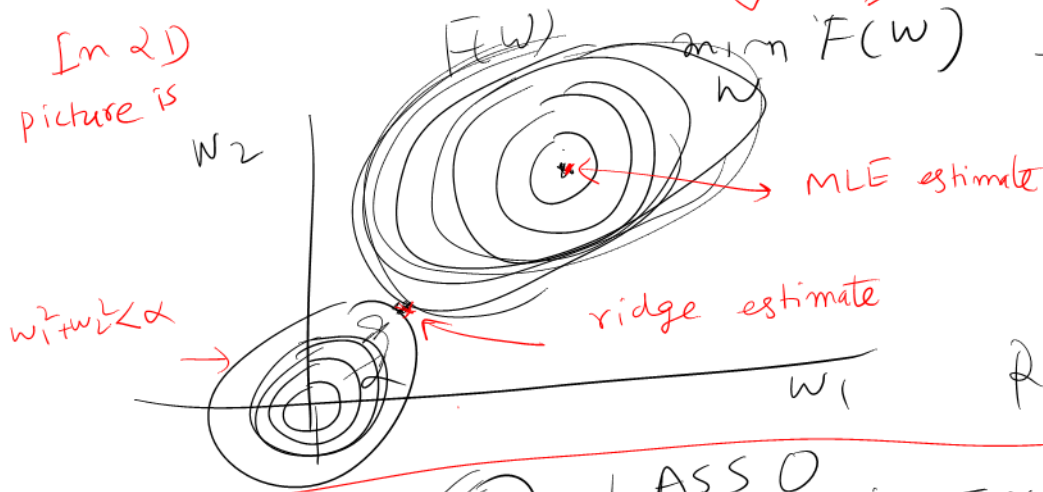we can force each $w_i$ with a prior belief i.e $w_i \sim N(x \mid 0, \tau)$


$N(x \mid 0, \tau)$

$$P(W \mid D)$$

then using MAP estimate of $W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix}$

equivalent to solving following optimization problem.

$$\arg\min_{W} \left( \underbrace{\sum_{i=1}^{N} (y_i - W^T x_i)}_{MLE} + \lambda \|W\|_2^2 \right)$$

$$\arg\min \left( \underbrace{\sum_{i=1}^{N} (y_i - W^T x_i)}_{F(W)} \right) \text{ s.t } \|W\|_2^2 < \alpha$$

$$\min_{W} F(W) \quad \text{s.t} \quad \|W\|_2 < \alpha$$

$$w_1^2 + w_2^2 < \alpha$$

[In 2D] picture is

$W_2$

$w_1^2 + w_2^2 < \alpha$

MLE estimate

ridge estimate

$W_1$

Ridge Regression

LASSO $\min F(W) + \lambda |W|$

MLE estimate

LASSO estimate of $W$

$W = |W_1| + |W_2|$
$\sim |W|_2$

$$\mu_1 \|W\|^2 + \mu_2 |W|$$

$\begin{bmatrix} W_1 \\ W_2 \end{bmatrix}$   $W$   $W_1$