$$w^* = \hat{w}_{ML} \quad \leftarrow \text{Max likelihood.} = \boxed{(XX^T)^\dagger Xy} \rightarrow \text{random}$$

$$y/x = \underline{w}^T x + \underline{\epsilon} \rightsquigarrow N(0, \boxed{\sigma^2}) \qquad \{(x_1, y_1), \cdots (x_n, y_n)\}$$

$$N(w^T x, \boxed{\sigma^2}) \qquad \qquad \underset{w^T x_i + \epsilon_i}{\downarrow}$$

$$\underline{w} \in \mathbb{R}^d \quad ; \hat{w}_{ML} \in \mathbb{R}^d$$

Want a way to understand how good $\hat{w}_{ML}$ is in estimating $w$

Mean Squared $\rightarrow \quad \mathbb{E}\left[\|\hat{w}_{ML} - w\|^2\right] = \sigma^2 \cdot \text{trace}\left((XX^T)^{-1}\right)$
error
$\qquad \qquad \qquad \underset{\text{over randomness in } y}{\downarrow}$

---

$$A = \begin{bmatrix} a_1 & & \\ & a_2 & \\ & & \ddots \\ & & & a_d \end{bmatrix}$$

$$tr(A) = \sum_{i=1}^{d} a_i = \sum_{i=1}^{d} \lambda_i$$

$\qquad \qquad \underset{\substack{i^{th} \text{ Eigenvalue} \\ \text{of } A}}{\downarrow}$

$$\boxed{\text{trace}\left((XX^T)^{-1}\right)}$$

Let Eigenvalues of $(XX^T)$ be $\{\lambda_1, \cdots, \lambda_d\}$

Eigenvalues of $(XX^T)^{-1}$ are $\left\{\frac{1}{\lambda_1}, \cdots, \frac{1}{\lambda_d}\right\}$

Mean sq. error $(\hat{w}_{ML})$

$$\mathbb{E}\left(\|\hat{w}_{ML} - w\|^2\right) = \sigma^2 \left(\sum_{i=1}^{d} \frac{1}{\lambda_i}\right)$$

---

Consider the following estimator:

$$\hat{w}_{new} = \left(XX^T + \lambda I\right)^{-1} Xy$$
$\qquad \qquad \qquad \underset{\in \mathbb{R}_+}{\downarrow} \quad \underset{\in \mathbb{R}^{d\times d}}{\searrow}$

$$\hat{w}_{ML} = (X\underline{X}^T)^{-1} Xy$$

$A v_i = \lambda_i v_i$

$(A + \lambda I)\underline{v}_i = A v_i + \lambda v_i$
$\qquad \qquad = \lambda_i v_i + \lambda v_i$
$\qquad \qquad = (\lambda_i + \lambda) \underline{v}_i$

For some matrix $A$, let Eigen values be $\{\lambda_1, \cdots, \lambda_d\}$

· What are Eigenvalues of $A + \lambda I$? $\{\lambda_1 + \lambda, \cdots, \lambda_d + \lambda\}$

$$\text{trace}\left((X\underline{X}^T + \lambda I)^{-1}\right) = \left(\sum_{i=1}^{n} \frac{1}{\lambda_i + \lambda}\right)$$

EXISTENCE Thm : (Informal)

$\exists \lambda \in \mathbb{R}$ s.t

$$\hat{w}_{new} = (XX^T + \lambda I)^{-1} Xy \quad \text{has lesser mean sq. error}$$

than $\hat{w}_{ML}$

In practice, find $\lambda$ by <u>CROSS VALIDATION</u>

|  | VALIDATION SET |
|---|---|
| TRAIN SET | |
| 80% | 20% |

- Train on the training set and check for error on validation set.
- Pick $\lambda$ that gives least error.

K-FOLD CROSS VALIDATION

| F1 | F2 | | ... | | Fk |
|---|---|---|---|---|---|

- Train on Folds $\{F_1, \dots F_{i-1}, F_{i+1} \dots F_k\}$
- Validate on $F_i$

- Pick $\lambda$ that gives least average error.

LEAVE ONE OUT CROSS VALIDATION

$$\hat{w}_{new} = (XX^T + \lambda I)^{-1} Xy$$

IS there an alternate way to understand $\hat{w}_{ML}$?

BAYESIAN MODELING

- NEED A PRIOR on $w$ i.e., $P(w)$ $\mathbb{R}^d$

$$\underline{\text{LIKELIHOOD}} \qquad y/x \sim N\left(w^T x, \boxed{1}\right)$$

$$\underline{A \quad \text{CHOICE} \quad \text{FOR} \quad \text{PRIOR}}$$

$$w \sim N\left(0, \gamma^2 I\right) \qquad \qquad \begin{bmatrix} \gamma^2 & & \\ & \gamma^2 & 0 \\ 0 & & \ddots \\ & & \gamma^2 \end{bmatrix}$$

$\in \mathbb{R}^d$

$$\gamma^2 I \longrightarrow \text{COVARIANCE} \quad \text{MATRIX} \quad \mathbb{R}^{d \times d}.$$

As usual,

$$P\left(w \mid \{(x_1, y_1) \cdots (x_n, y_n)\}\right) \propto P\left(\{(x_1, y_1), \cdots, (x_n, y_n)\} / w\right) \cdot P(w)$$

$$\propto \left(\prod_{i=1}^{n} e^{-\frac{(y_i - w^T x_i)^2}{2}}\right) \cdot \left(\prod_{i=1}^{d} e^{-\frac{(w_i - 0)^2}{2\gamma^2}}\right)$$

$$e^{-\sum_{i=1}^{d} \frac{w_i^2}{2\gamma^2}}$$

$$\propto \left(\prod_{i=1}^{n} e^{-\frac{(y_i - w^T x_i)^2}{2}}\right) \cdot e^{-\frac{\|w\|^2}{2\gamma^2}}$$

How will the MAP estimate look like?

$$\hat{w}_{MAP} = \underset{w}{\arg\max} \sum_{i=1}^{n} -\frac{(y_i - w^T x_i)^2}{2} - \frac{\|w\|^2}{2\gamma^2}$$

$$\hat{w}_{MAP} = \underset{w}{\arg\min} \frac{1}{2} \sum_{i=1}^{n} (y_i - w^T x_i)^2 + \frac{1}{2\gamma^2} \|w\|^2 \longrightarrow f(w)$$

Take gradient, set it to 0 to Solve for $\hat{w}_{MAP}$.

$$\nabla f(\omega) = (XX^T)\omega - Xy + \frac{\omega}{\gamma^2} \qquad \boxed{\text{verify}}$$

$$\boxed{\hat{\omega}_{MAP} = \left(XX^T + \underbrace{\frac{1}{\gamma^2}}_{} I\right)^{-1} Xy}$$

CROSS VALIDATE in PRACTICE.

CONCLUSION: MAP ESTIMATION for linear regression with a Gaussian prior $\boxed{N(0, \gamma^2 I)}$ for $\omega$ is Equivalent to "NEW" estimator we used earlier.

---

LINEAR REGRESSION

$$\hat{W}_{ML} = \underset{\omega}{\arg\min} \; \sum_{i=1}^{n} \left(\hat{W}^T z_i - y_i\right)^2$$

$\boxed{\frac{1}{\gamma^2}}$

RIDGE REGRESSION

$$\hat{W}_R = \underset{\omega}{\arg\min} \; \underbrace{\sum_{i=1}^{n} \left(\hat{W}^T z_i - y_i\right)^2}_{LOSS} + \underbrace{\lambda \|W\|^2}_{REGULARIZER}$$

RIDGE

| $f_1$ height | $f_2$ weight | $f_3$ 2 height + 3 weight | label 3 height+4 weight |
|---|---|---|---|
| ① | ① | ① | |
| 0 | $c_1$ | $c_2$ | |
| 3 | 4 | ⓪ | |