

DATASET : $\{(x_i, t_i)\}_{i=1}^N$, $x_i \in \mathbb{R}$ = samples
 $t_i \in \mathbb{R}$ = labels

FEATURES

Polynomial
basis fct.

$$\phi: \mathbb{R} \rightarrow \mathbb{R}^{M+1}$$

$$x \mapsto \begin{bmatrix} x^0 \\ x^1 \\ \vdots \\ x^M \end{bmatrix}$$

MAPPING fct : $y(x_i) = \sum_{j=0}^M w_j \cdot x_i^j$, f_i

$$y = X \cdot w$$

FEATURE
MATRIX : \tilde{X} is $N \times (M+1)$

$$\tilde{X} = \begin{bmatrix} x_1^0 & x_1^1 & x_1^2 & \dots & x_1^M \\ x_2^0 & x_2^1 & x_2^2 & \dots & x_2^M \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_N^0 & x_N^1 & x_N^2 & \dots & x_N^M \end{bmatrix}$$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_M \end{bmatrix}$$

OBJECTIVE : $J(\omega) = \frac{1}{2} \sum_{i=1}^N \epsilon_i^2$

for

where $\epsilon_i = t_i - y(x_i) = t_i - y_i$

$$J(\omega) = \frac{1}{2} \| t - X \cdot \omega \|_2^2$$

t is the labels vector, $N \times 1$, $t = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix}$

Solution
for
parameters
 ω

\bar{X} is $N \times (\mu+1)$

$$\frac{\partial J(\omega)}{\partial \omega} = 0$$
$$\Leftrightarrow \boxed{\omega = (\bar{X}^T \bar{X})^{-1} \bar{X}^T t} \quad (\mu+1) \times 1$$

↑
Assuming the inverse
of $\bar{X}^T \bar{X}$ exists

PSEUDO-CODE for LINEAR REGRESSION:

INPUT: input DATA $\{x_i\}_{i=1}^N = \mathbf{x}$
target values $\{t_i\}_{i=1}^N = \mathbf{t}$
hyperparameter $\lambda \in \mathbb{R}$

TRAINING
FUNCTION

- ① Compute the feature matrix \mathbf{X}
- ② Compute the ~~solution~~ ^{parametric} for w :
 $w = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$
- ③ Predict labels for samples x
 $y = \mathbf{X} \cdot w$
- ④ Compute the residual error
 $e = t - y$

Return w, y, e

During test ...

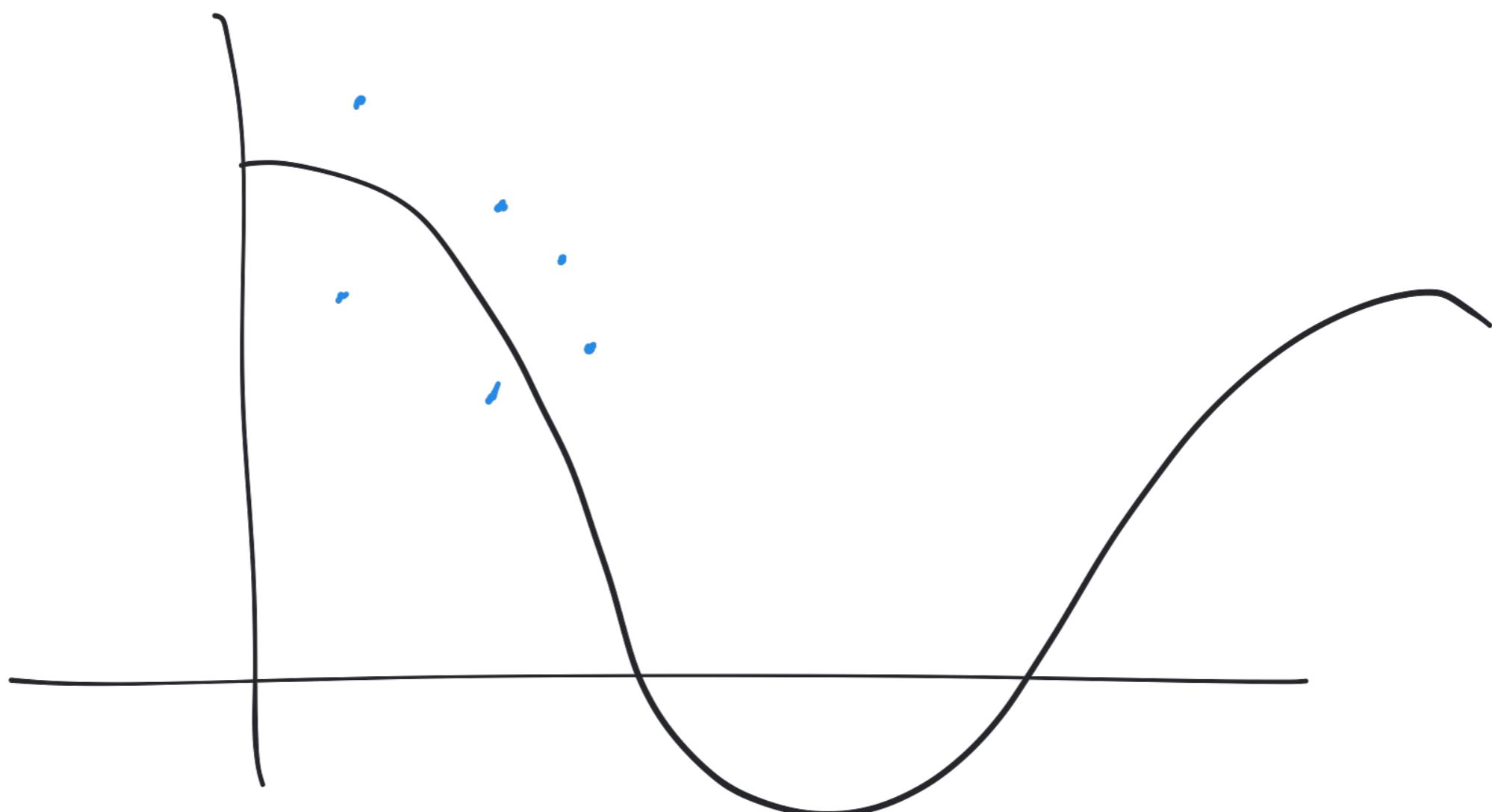
INPUT: TEST samples $\{z_i\}_{i=1}^N = z$

(trained) parameters w .

① Construct feature matrix for
test samples

$$Z_1 = \begin{bmatrix} z_1^0 & z_1^1 & \dots & z_1^M \\ z_2^0 & z_2^1 & \dots & z_2^M \\ \vdots & \vdots & \ddots & \vdots \\ z_N^0 & z_N^1 & \dots & z_N^M \end{bmatrix}$$

② $y_{TEST} = Z_1 \cdot w$



OBJECTIVE function:

will have 2 terms

error term

$$J(\omega) = \frac{1}{2} \|t - X\omega\|_2^2$$

+ regularization term

$$\sum_{j=0}^N \omega_j^2$$

hyperparameter

ridge
regularizer

$\lambda \rightarrow 0$: no penalty on the parameters ω

$\lambda \rightarrow \infty$: forces the parameters +
be small.

Lasso regularizer: $R(w) = \sum_{j=0}^M |w_j|$
 or
 (L1-penalty)

Ridge regularizer: $R(w) = \sum_{j=0}^M w_j^2$
 or
 (L2-penalty)

Elastic Net: $R(w) = \beta \sum_{j=0}^M |w_j| + (1-\beta) \cdot \sum_{j=0}^M w_j^2$

β is an additional hyperparameter
 for Elastic Net

Regularizer is also known as a
 penalty term.

Observations:

Lasso: ① Drive the parameters ^{to}
exactly 0 much faster than ridge

② promotes sparsity
↳ a form of feature selection

Ridge:

① will apply a stronger penalty to
outliers compared to Lasso.
② will "force" parameters to be
smaller but never 0.

$$\textcircled{1} \quad \omega = [0.5, 0.5, 1]$$

$$\|\omega\|_2^2 \equiv \text{ridge} = 0.5^2 + 0.5^2 + 1^2 = 1.5$$

$$\|\omega\|_1 \equiv \text{Lasso} = |0.5| + |0.5| + |1| = 2$$

$$\textcircled{2} \quad \omega = [0, 0, 2]$$

$$\|\omega\|_2^2 = 0^2 + 0^2 + 2^2 = 4 \equiv \text{ridge}$$

$$\|\omega_1\|_1 = |0| + |0| + |2| = 2 \equiv \text{Lasso}$$

$$y = w_0 + w_2 x^2$$

Objective function w/ Ridge Regularizer

$$J(w) = \frac{1}{2} \|t - Xw\|_2^2 + \frac{\lambda}{2} \|w\|_2^2$$

optimization
problem

$$\arg \min_w J(w)$$

$$J(\omega) = \frac{1}{2} (t - \mathbf{x}\omega)^T(t - \mathbf{x}\omega) + \frac{\lambda}{2} \omega^T\omega$$

$$= \frac{1}{2} (t^T - \omega^T \mathbf{x}^T)(t - \mathbf{x}\omega) + \frac{\lambda}{2} \omega^T\omega$$

$$= \frac{1}{2} (t^T t - t^T \mathbf{x}\omega - \omega^T \mathbf{x}^T t + \omega^T \mathbf{x}^T \mathbf{x}\omega) + \frac{\lambda}{2} \omega^T\omega$$

— II —

$$\begin{aligned} \frac{\partial J(\omega)}{\partial \omega} &= 0 \Leftrightarrow \frac{1}{2} \left(-t^T \mathbf{x} - (\mathbf{x}^T \cdot t)^T + (\mathbf{x}^T \mathbf{x}\omega)^T \right. \\ &\quad \left. + \omega^T \mathbf{x}^T \mathbf{x} \right) + \frac{\lambda}{2} \cdot \omega^T + \frac{\lambda}{2} \cdot \omega^T = 0 \end{aligned}$$

$$\Leftrightarrow -t^T \cdot \mathbf{x} + \omega^T \mathbf{x}^T \mathbf{x} + \lambda \cdot \omega^T = 0$$

$$\Leftrightarrow \omega^T \mathbf{x}^T \mathbf{x} + \lambda \omega^T = t^T \cdot \mathbf{x}$$

$$\Leftrightarrow \cancel{X^T X} w + \lambda \cdot w = \cancel{X^T} \cdot t$$

applying
transpose
on both
sides

$$\Leftrightarrow (X^T X + \lambda \cdot I) \cdot w = \cancel{X^T} \cdot t$$

if inverse exists

identity matrix
 $(n+1) \times (n+1)$

$$\Leftrightarrow w = (X^T X + \lambda \cdot I)^{-1} \cdot \cancel{X^T} \cdot t$$

$$w = (X^T X)^{-1} X^T t$$

If $(X^T X)^{-1}$ does not exist

① $\det(X^T X) = 0$

② $X^T X$ is NOT full rank

③ $X^T X$ has linearly dependent columns

④ This will happen when $M > N$.

⑤ Solution: diagonally-load $X^T X$

before inverting it

↳ this is the solution for
least squares w/ RIDGE REG.

PSEUDO-CODE:

INPUT: samples $x = \{x_i\}_{i=1}^N$,
target $t = \{t_i\}_{i=1}^N$
hyperparameters N and λ

- ① Compute feature matrix $\underline{\mathbf{X}}$
- ② Compute solution for w :
 $w = (\underline{\mathbf{X}}^T \underline{\mathbf{X}} + \lambda \cdot I)^{-1} \cdot \underline{\mathbf{X}}^T \cdot t$
- ③ Model prediction:
 $y = \underline{\mathbf{X}} \cdot w$

Retorna w, y

Ways to avoid over fitting

- ① Apply Occam's Razor principle
- ② Add more data
- ③ Use regularization term
- ④ Cross-validation

