# Regularization

## Ridge

$$R_w(\omega) = \sum_i^+ w_i^2$$

## Lasso

$$R_w(\omega) = \sum_i^+ |w_i|$$

$$J(\omega) = \boxed{\frac{1}{2} \sum_n^1 e_n^2} + \boxed{\frac{1}{2} \lambda \cdot R_w(\omega)}$$

$$\boxed{\underset{w}{\arg\min} \; J(\omega)}$$

Regularization weight

Ridge : ① Will penalize large
weight values more

Lasso

② Highly affected
by outliers

③ prefers that coefficients
are not zero
but very small

$R(w)$

Error ridge

Error
Lasso

Ridge

Lasso

W

Lasso: ① Less sensitive to outliers

② Makes weights go to zero much
faster

③ prefers some coefficients exactly
zero.

Using Lasso reg. we can perform Feature Selection

$$J(w) = J_E(x,w) + \lambda \cdot R_w(w)$$

$\lambda = 0$ : only minimizing $J_E(x,w)$

$\lambda \to \infty$ : disregards $J_E(x,w)$ and forces $R_w$ to be small

FOR Regularized Polynomial Regression:

we control: ① Model
② Cost function
③ Learning algorithm

Model order $\underline{\underline{M}}$ and reg. weight $\underline{\underline{\lambda}}$

$$J(w) = \frac{1}{2} \sum_{i=1}^{N} \underbrace{(t_i - y_i)^2}_{= \ell_i} + \frac{1}{2} \cdot \lambda \cdot \sum_{j=0}^{M} w_j^2$$

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

If we have $N \gg M$ : likely will

# samples    model order    also be an overdetermined system of eqs

→ identity matrix

$$\lambda \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & & \ddots & 1 \\ 0 & & & 1 \end{bmatrix} \begin{array}{c} (M+1) \times \\ (M+1) \end{array}$$

$$w^* = \left( X^T X + \lambda \cdot I \right)^{-1} \cdot X^T \cdot t$$

REGULARIZED polynomial REGRESSION

$$w^* = \left( X^T X \right)^{-1} X^T \cdot t$$

polynomial reg. w/out reg.

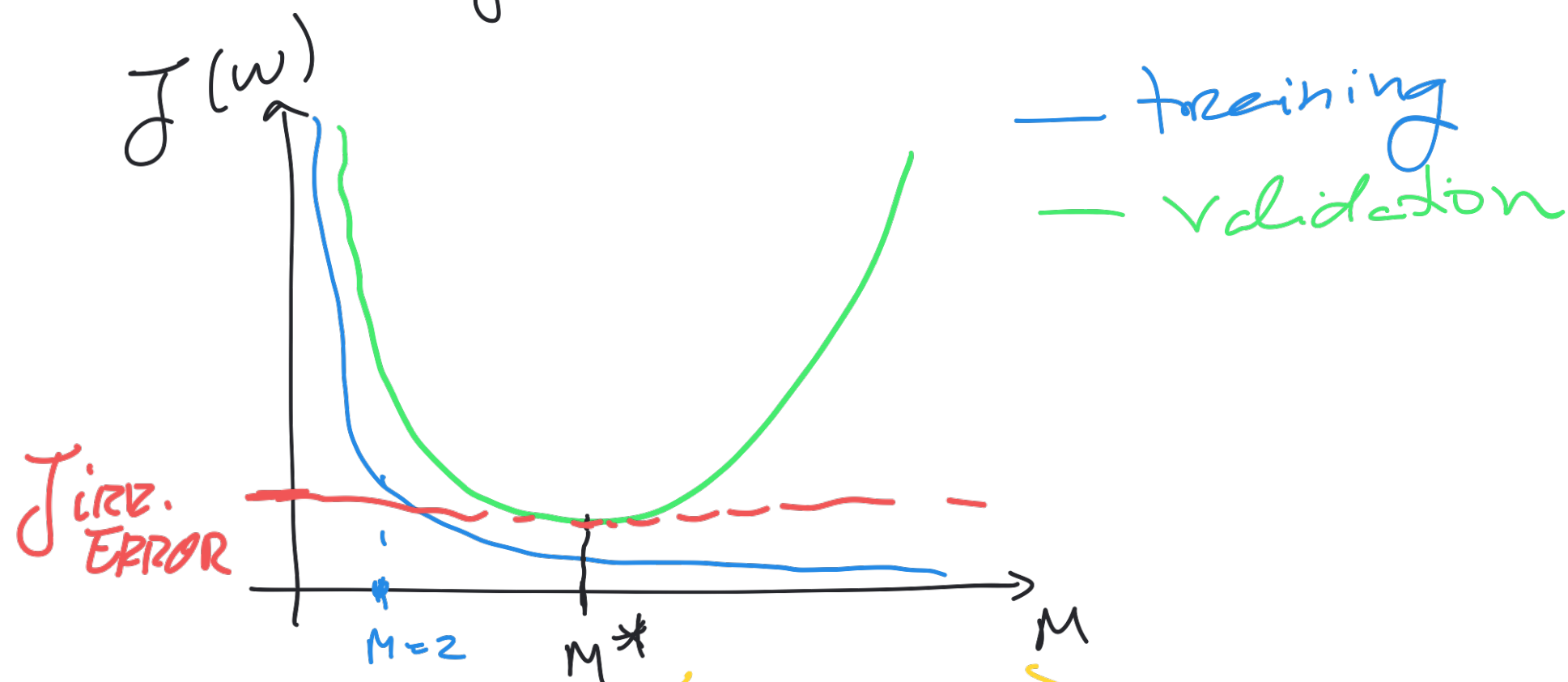Regularization:

$\rightarrow$ diagonally loading $X^T X$

$\rightarrow$ make such $X^T X$ is full rank

$$w_0 \cdot G + w_1 \cdot G + w_2 G + w_3 G$$

$$y = \sum_{j=0}^{M-1} w_j \cdot \phi_j(x)$$

- We are looking for the best combination of controllable parameters $(M, \lambda)$:

$J(w)$

— training
— validation

$J_{irr.\ error}$

M=2    M*    M

← underfitting    overfitting →

$J_{train}$ will be high

$J_{val}$ will be high

high bias
low variance

$J_{train}$ will be small

$J_{val} \gg J_{train}$

low bias &
high varian