

Regularization

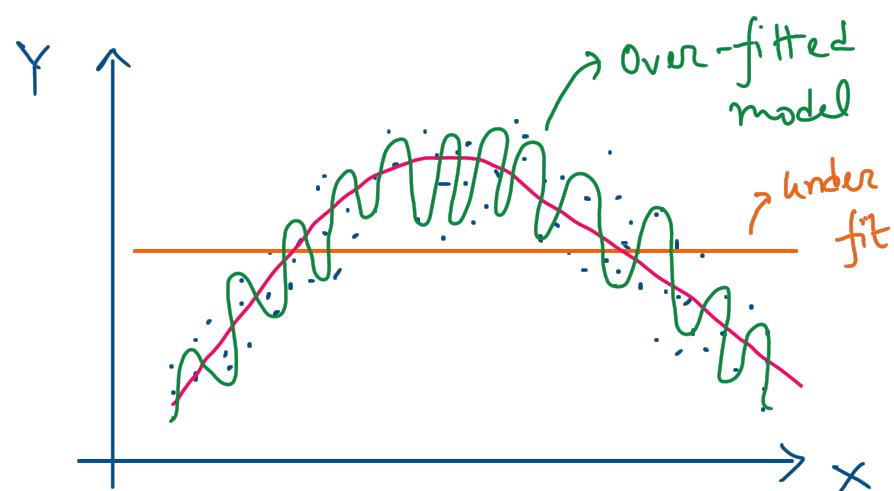
What is regularization?

Regularization is a technique to reduce the overfitting of a model.

- Overfitted models fails to generalize → the model performs well in training dataset, but doesn't perform well in test dataset.
- Regularization tries to reduce the overfitting by reducing the variance in the dataset.

What is the cause of overfitting?

- Less number of training datapoints.
- Too many features
- Too complex model.



Linear Regression

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_k x_k$$

$\frac{\Delta \hat{y}}{\Delta x_j} = \theta_j$ } interpretation of model weights : change in response per unit change in feature.

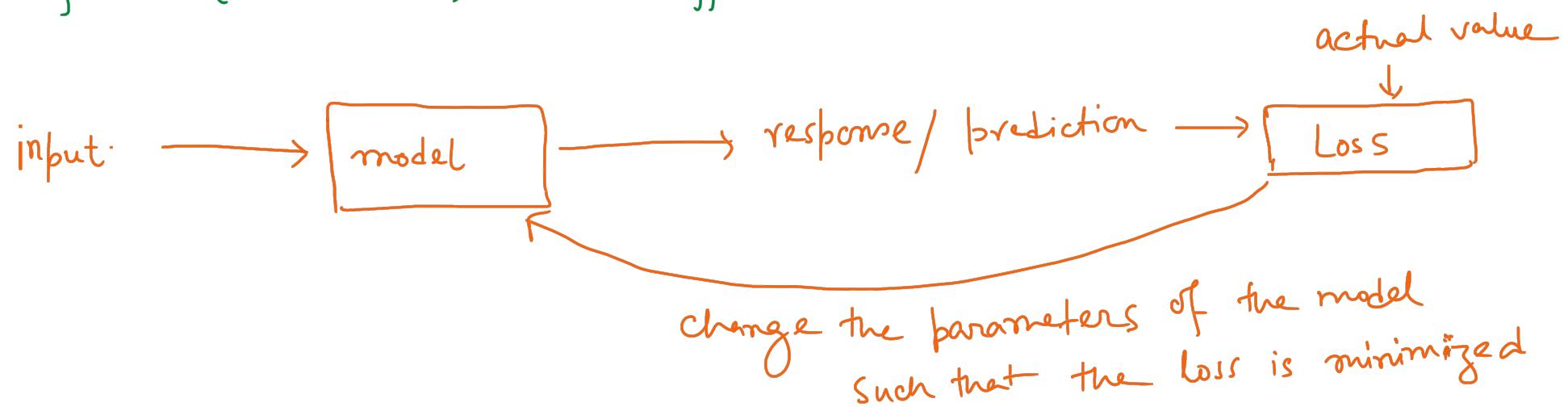
if $\underline{\theta_j}$ is high \rightarrow model response changes is high when the corresponding feature changes. | model is more susceptible to feature x_j

θ_j is small \rightarrow , ^{change in} model response is low when the corresponding feature changes. | model is less susceptible to feature x_j

Logistic regression: $\hat{y} = \text{sigmoid}(z)$ where $z = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_k x_k$

$$\theta_j : \frac{\Delta z}{\Delta x_j}$$

Regularization : Is a way to limit the parameters (in case of a parametric model like linear or logistic regression) so that slight change in features (due to noise) doesn't affect the model response much.



We introduce a "penalty" term in the loss-function. When the model parameters are large the penalty will be large & when the model parameter is low the penalty will be low.

$$\text{"penalty"} = f(\text{model weights / parameters})$$

$$\text{Loss}(\Theta) = \frac{1}{m} \sum_{i=1}^m \text{error}(\hat{y}^{(i)}, y^{(i)})$$

for linear regression: $\text{Loss}(\underline{\Theta}) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$

for logistic regression (binary): $\text{Loss}(\Theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})]$

for softmax regression: $\text{Loss}(\Theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^k y_j^{(i)} \log(\hat{y}_j^{(i)})$

Regularized Loss function: $L(\underline{\Theta}, \alpha) = \text{Loss}(\underline{\Theta}) + \boxed{\alpha \cdot (\text{penalty})}$

regularization
parameter

Regularization term.

L2 - regularization :

$$L(\theta, \alpha) = \frac{1}{m} \sum_{i=1}^m \text{error}(\hat{y}^{(i)}, y^{(i)}) + \alpha \left(\frac{1}{2} \sum_{j=1}^k \theta_j^2 \right)$$

penalty term : $\frac{1}{2} \sum_{j=1}^k \theta_j^2 = \frac{1}{2} (\theta_1^2 + \theta_2^2 + \theta_3^2 + \dots + \theta_k^2)$

We don't use bias in the penalty : In the penalty term we use only those parameters which are associated with the inputs of the model.

Regularization parameter (α) : How much penalty is injected to the regularized - loss function. (usually $\alpha > 0$)

if $\alpha = 0$: no-penalty (unregularized model)

α is high : high - penalty

L2 Regularized MSE loss:

$$L(\Theta, \alpha) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 + \frac{\alpha}{2} \sum_{j=1}^k \theta_j^2$$

Ridge-regression (L2-regularized linear regression)

$$\frac{\partial L}{\partial \theta_j} = \begin{cases} \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) & ; \text{for bias } \theta_0 \\ \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) x_j^{(i)} + \alpha \theta_j & ; \text{for weight } \theta_j \end{cases}$$

$$\begin{aligned} \theta_j(t+1) &\leftarrow \theta_j(t) - \eta \frac{\partial L}{\partial \theta_j} = \theta_j(t) - \eta \left[\frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) x_j^{(i)} + \alpha \theta_j(t) \right] \\ &= (1 - \eta \alpha) \theta_j(t) - \frac{\eta}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) x_j^{(i)} \end{aligned}$$

L1-regularization

$$\vec{v} = [v_1, v_2, v_3, \dots, v_k]^T$$

$$L2\text{-norm} = \|\vec{v}\|_2 = \sqrt{v_1^2 + v_2^2 + \dots + v_k^2}$$

$$L1\text{-norm} = \|\vec{v}\|_1 = |v_1| + |v_2| + \dots + |v_k|$$

$$L(\theta, \alpha) = \frac{1}{m} \sum_{i=1}^m \text{error}(\hat{y}^{(i)}, y^{(i)}) + \alpha \sum_{j=1}^k |\theta_j|$$

L1 regularized MSE Loss : $L(\theta, \alpha) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 + \alpha \sum_{j=1}^k |\theta_j|$

LASSO Regression

\downarrow Least Absolute Shrinkage and Selection Operator.

feature selection

Combination of L_1 & L_2 :

$$L(\theta, \alpha, \beta) = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 + \alpha \left[\underbrace{\beta \sum_{j=1}^k |\theta_j|}_{L_1} + \underbrace{\frac{1-\beta}{2} \sum_{j=1}^k \theta_j^2}_{L_2} \right]$$

$\alpha \rightarrow$ regularization parameter ($\alpha=0$ means no regularization)

$\beta \rightarrow L_1$ ratio \rightarrow it controls the presence of L_1 & L_2 regularization in the penalty

$0 \leq \beta \leq 1 \rightarrow \beta=0 : \text{no } L_1, \text{ only } L_2$

$\beta=1 : \text{only } L_1, \text{ no } L_2$

} Elastic - net

Regularized MSE cost:

$$L(\theta, \alpha, \beta) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 + \alpha \left[\beta \sum_{j=1}^k |\theta_j| + \frac{1-\beta}{2} \sum_{j=1}^k \theta_j^2 \right]$$

} Regression