Linear Regression: Loss function

Fitting an estimator/predictor/model involves solving for the Θ that minimizes the Loss function.

Recall our goal is to make the discrepancy (error) between $\hat{\mathbf{y}}$ and $\hat{\mathbf{y}}$ "small".

• The discrepancy between $\mathbf{y^{(i)}}$ and $\hat{\mathbf{y}^{(i)}}$ is referred to as the *residual*, usually denoted by ϵ

$$\epsilon^{(\mathbf{i})} = \mathbf{y^{(i)}} - \hat{\mathbf{y}}^{(\mathbf{i})}$$

So

$$\mathbf{y} = \hat{\mathbf{y}} + \epsilon$$
 $= \mathbf{X}\Theta + \epsilon$

We define the per-example loss to be the residual squared

$$\mathcal{L}_{\Theta}^{(\mathbf{i})} = (\mathbf{y^{(i)}} - \hat{\mathbf{y}^{(i)}})^2$$

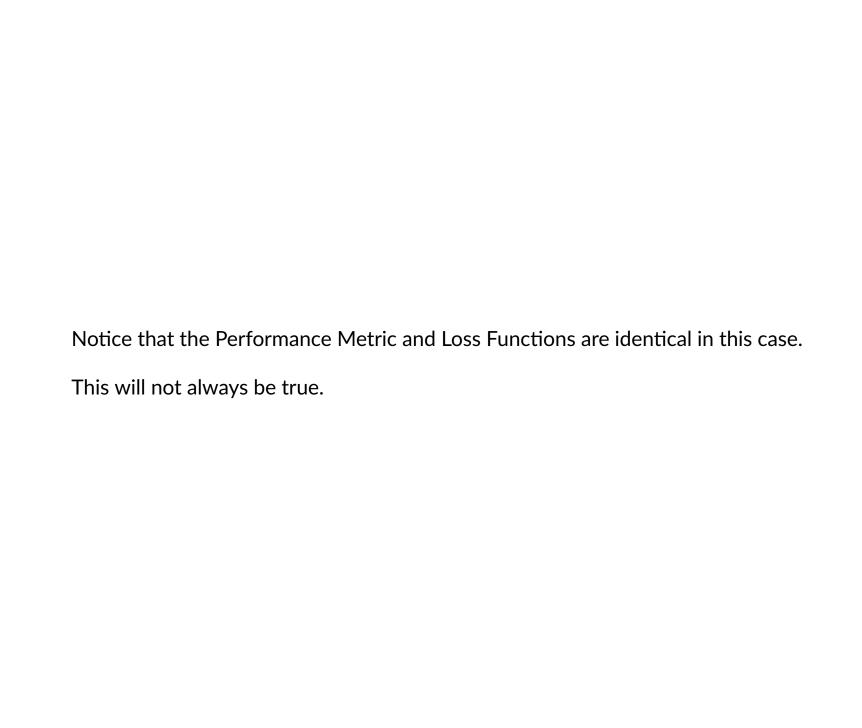
so that the average loss

$$egin{array}{lll} \mathcal{L}_{\Theta} & = & rac{1}{m} \sum_{i=1}^m \mathcal{L}_{\Theta}^{(\mathbf{i})} \ & = & rac{1}{m} \sum_{i=1}^m (\mathbf{y^{(i)}} - \hat{\mathbf{y}^{(i)}})^2 \end{array}$$

This expression on the right is called the *Mean Squared Error (MSE)*.

$$\mathrm{MSE}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{m} \sum_{i=1}^{m} (\mathbf{y^{(i)}} - \hat{\mathbf{y}^{(i)}})^2$$

• You will sometimes see *Root Mean Squared Error (RMSE)* which is the square root of the MSE



$oldsymbol{R^2}$ versus RMSE: Absolute versus relative error

One often sees the term \mathbb{R}^2 in the context of Linear Regression.

Whereas RMSE is absolute error (in same units as \mathbf{y}), R^2 is a relative error (in units of percent).

The relationship is:

$$egin{array}{lll} R^2 & = & 1 - \left(rac{\sum_{i=1}^m \left(y_i - \hat{y}_i
ight)^2}{\sum_{i=1}^m \left(y_i - ar{y}_i
ight)^2}
ight) \ & = & 1 - \left(rac{m \cdot ext{MSE}(\mathbf{y}, \hat{\mathbf{y}})}{\sum_{i=1}^m \left(y_i - ar{y}_i
ight)^2}
ight) \ & = & 1 - \left(rac{m \cdot ext{RMSE}(\hat{y}, y)^2}{\sum_{i=1}^m \left(y_i - ar{y}_i
ight)^2}
ight) \end{array}$$

In addition to changing the units of error, the ${\cal R}^2$ metric has an interesting interpretation.

Consider a naive "baseline" model for prediction

- predict $\bar{\mathbf{y}}$ for every value of \mathbf{x}
 - where $\bar{\mathbf{y}}$ is the average (over the training examples) of the target

The loss for the naive model is

$$\mathcal{L}_{ ext{naive}} = ext{MSE}(\mathbf{y}, ar{\mathbf{y}})$$

Then

$$egin{array}{lll} R^2 &=& 1 - \left(rac{m \cdot ext{MSE}(\mathbf{y}, \hat{\mathbf{y}})}{m \cdot ext{MSE}(\mathbf{y}, ar{\mathbf{y}})}
ight) \ &=& 1 - rac{\mathcal{L}}{\mathcal{L}_{ ext{paive}}} \end{array}$$

Thus, R^2 is the percent reduction in loss achieved by our model compared to the naive model that always predicts $\bar{\mathbf{y}}$.

We now know our Loss function.

The "solution" to the Linear Regression task is finding ("fitting") the Θ that minimizes average loss

$$\Theta = \operatorname*{argmin}_{\Theta} \mathcal{L}_{\Theta}$$

which is the Θ that minimizes the MSE.