

The Learning Penalty
vs.

Measuring Prediction Performance

Similar in form but
Different in function.

The Learning Penalty (aka the Penalty function)

			$w_0 x_0 + w_1 x_1 + w_2 x_2$	Penalty
	Feature1	Feature2	Output	$(y - \hat{y})^2$
row 1	x_1	x_2	y	\hat{y}
:	:	:	:	:
row m				

What you deal with
in orange

Cost = Sum of all penalties

(add up the value of
each row)

More precisely, Cost = $\frac{1}{m} (\text{sum of all penalties})$

this makes it the average penalty

How Learning Happens

Iteration 1

usually just randomly chosen values

$$w_0 = 0.5, w_1 = 2.2, w_2 = -1.8$$



Penalty $(y - \hat{y})^2$ for each row

Add up the values

Divide by M (# of rows)



Cost of iteration

{ Goal: Find the next set of w_0, w_1, w_2 that lower this cost.

Gradient descent is the method for doing this.

How Learning Happens (Contd.)

Iteration 2

$$w_0 = 0.52, w_1 = 2.1, w_2 = -1.92$$

→ Penalty $(y - \hat{y})^2$ for each row

Add up the values

Divide by M (# of rows)



Cost of iteration

Goal: find the next set of
 w_0, w_1, w_2 that lower
this cost.



Iteration 3, 4, 5, 6, ..., 23

We can choose to stop after a certain #
of iterations, or when cost isn't
going down by a lot any more.

Summary - For each iteration,

$$\text{Penalty} = (y - \hat{y})^2 \text{ for each row}$$

$$\text{Cost} = \frac{1}{m} \sum [(y - \hat{y})^2]$$

#rows
in the dataset

Average or Mean

Penalty per iteration

"Mean Squared" Penalty

Once we get to the last iteration,

we have the optimal values for
 w_0 , w_1 , and w_2 .

Prediction

what

I give you feature values x_1, x_2 .

You tell me what the output is.

How

Build a model { Take the feature values x_1, x_2 .
Get the optimal values of w_0, w_1, w_2 .

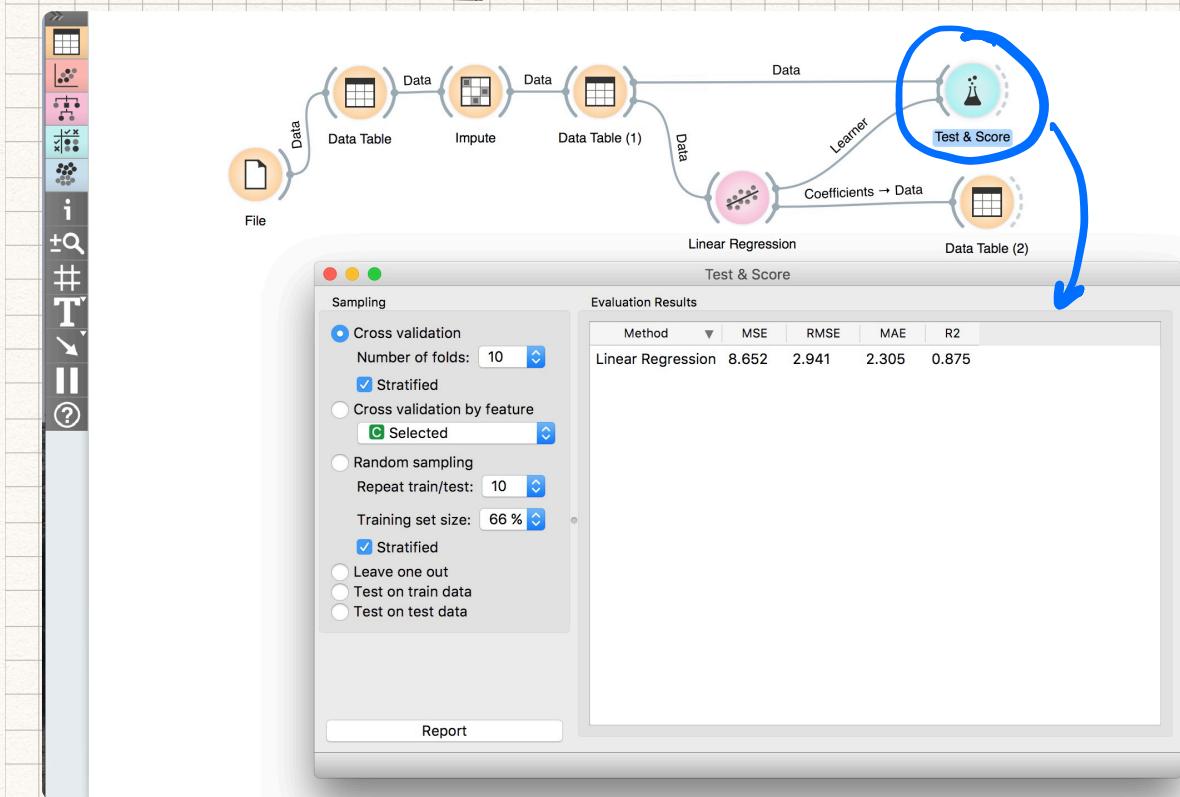
Predicted Output = $(w_0 * 1) + (w_1 * x_1) + (w_2 * x_2)$

always = 1

Use the model

Prediction (contd.)

How good is the prediction?



MSE = Mean Square Error

$RMSE$ = Root Mean Square Error

MAE = Mean Absolute Error

R^2 = R^2 = "R-Squared"

Always a good idea to have
a single number (if at all possible)
to measure how good your
predictions are.

MSE, RMSE, MAE and R^2 are
such numbers. Any one of
them can be used to measure
how good the predictions are.

The lower the value, the
better for MSE, RMSE, & MAE.

MSE - Mean Square Error

Features

x_1	x_2	Prediction	Actual	(Prediction - Actual)
				23.42
				5.8

→ This is
called the
"test" dataset

$$MSE = \frac{[(23.42 + 5.8)]}{2}$$

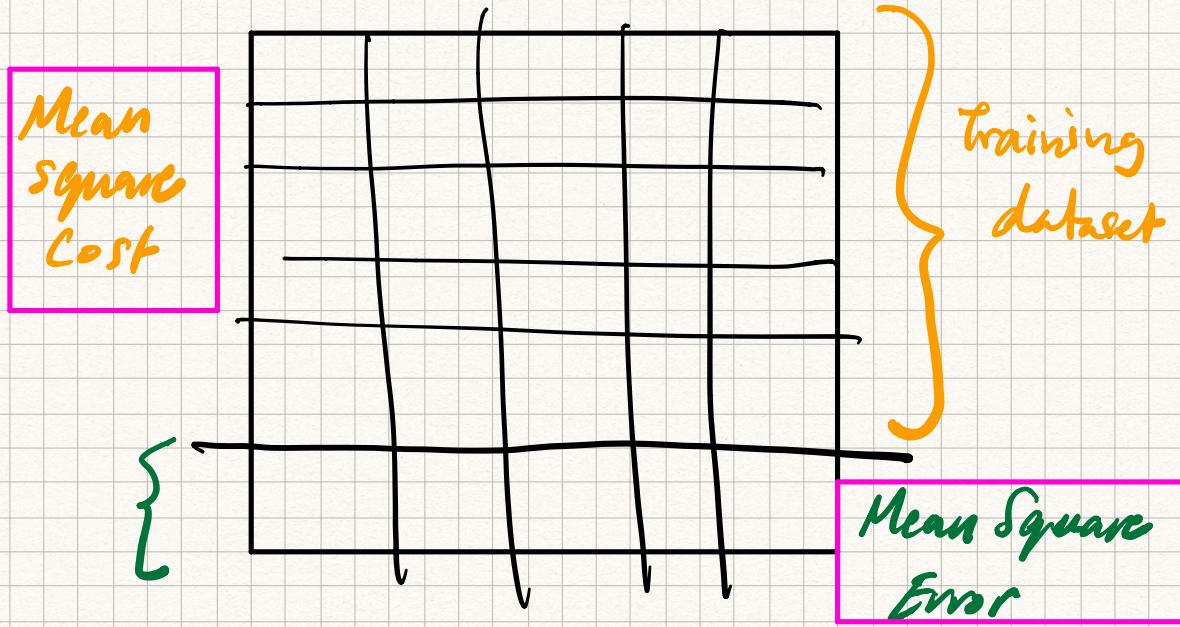


of rows

Identical to the Cost used in
gradient descent. But serves
an entirely different purpose.

Prediction Performance

Dataset



Test Dataset

- Use the **Training dataset** to calculate the optimal values of w_0 , w_1 , and w_2 .
- Use the **test dataset** to measure the prediction performance of the model.

SOME FORMULAS

(OPTIONAL)

SS = Sum of Squares

\bar{y} = average value of actual outputs

\hat{y} = predicted value of output

$y^{(i)}$ = actual output value
of the i th row of the
dataset.

$\hat{y}^{(i)}$ = predicted output value
of the i th row of the
dataset.

$|a - b|$ = Absolute value of
 $a - b$.

$$SS_{\text{total}} = \sum_{\text{all rows}} \left[(y^{(i)} - \bar{y})^2 \right]$$

$$SS_{\text{regression}} = \sum_{\text{all rows}} \left[(\hat{y}^{(i)} - \bar{y})^2 \right]$$

$$SS_{\text{residual}} = \sum_{\text{all rows}} \left[(y^{(i)} - \hat{y}^{(i)})^2 \right]$$

$$R^2 = 1 - \frac{SS_{\text{residual}}}{SS_{\text{total}}}$$

MSE = Mean Square Error

$$= \frac{SS_{\text{residual}}}{\# \text{ of rows in dataset}}$$

$RMSE$ = Root Mean Square Error

$$= \sqrt{MSE}$$

MAE = Mean Absolute Error

$$\frac{\text{Sum over all rows} \left[|\hat{y}^{(i)} - y^{(i)}| \right]}{\# \text{ of rows in the dataset}}$$

Absolute value