

MM 225 – AI and Data Science

Day 28: Logistic Regression 2

Instructors: Hina Gokhale, MP Gururajan, N. Vishwanathan

10 OCTOBER 2024

A solid blue horizontal bar spanning the width of the slide at the bottom.

Good news!

LSE function for multiple regression is a convex function:

$$SSE = \sum (y^{(i)} - \mathbf{w} \cdot \mathbf{x}^{(i)})^2$$

Cross entropy error function is also convex

$$\mathcal{E}(\mathbf{w}) = - \sum y^{(i)} \log(\psi(\mathbf{w} \cdot \mathbf{x}^{(i)})) + (1 - y^{(i)}) \log(1 - \psi(\mathbf{w} \cdot \mathbf{x}^{(i)}))$$

$$\psi(z) = \frac{1}{1 + e^{-z}}$$

Regression

$$SSE = \sum (y^{(i)} - \mathbf{w} \cdot \mathbf{x}^{(i)})^2 = [\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2(\mathbf{X}^T \mathbf{y})^T \mathbf{w}] + const$$

Initialize with \mathbf{w}_1 at random

set $t = 1$ and choose λ

continue until convergence

1. compute the gradient $\nabla SSE = \mathbf{X}^T (\mathbf{X} \mathbf{w}_t - \mathbf{y})$
2. update $\mathbf{w}_{t+1} = \mathbf{w}_t - \lambda \nabla SSE$
3. $t \leftarrow t + 1$

Logistic Regression

$$\mathcal{E}(\mathbf{w}) = - \sum y^{(i)} \log(\psi(\mathbf{w} \cdot \mathbf{x}^{(i)})) + (1 - y^{(i)}) \log(1 - \psi(\mathbf{w} \cdot \mathbf{x}^{(i)}))$$

$$\psi(z) = \frac{1}{1+e^{-z}} \text{ then } \frac{d\psi(z)}{dz} = \psi(z)(1 - \psi(z))$$

It can be shown that

$$\frac{\partial \mathcal{E}(\mathbf{w})}{\partial \mathbf{w}} = \sum_i [\psi(\mathbf{w} \cdot \mathbf{x}^{(i)}) - y^{(i)}] \mathbf{x}^{(i)}$$

Algorithm for Logistic Regression

1. Choose λ
2. Initiate \mathbf{w}_1
3. Iterate as $\mathbf{w}_{t+1} = \mathbf{w}_t - \lambda \sum_i [\psi(\mathbf{w}_t \cdot \mathbf{x}^{(i)}) - y^{(i)}] \mathbf{x}^{(i)}$
4. Iterate until convergence

Choosing λ in practice

Try out $\lambda = 0.001, 0.01, 0.1$ on a test data set. Choose the λ that gives stable and fast convergence.

To reach true minimum, reduce λ by factor of 10 as learning saturates.

Validation

Validation with test data

Model estimation using training data

Apply the same model to test data and see if the results are “same”!

What is meant by “same”?

Misclassification Quantification

Total data size = P + N		Predicted Condition	
		Positive(PP)	Negative (PN)
Actual Condition	Positive (P)	True Positive (TP)	False Negative (FN)
	Negative (N)	False Positive (FP)	True Negative (TN)

$$\text{TPR} = \text{True Positive Rate} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \text{Sensitivity / Hit Rate / Recall}$$

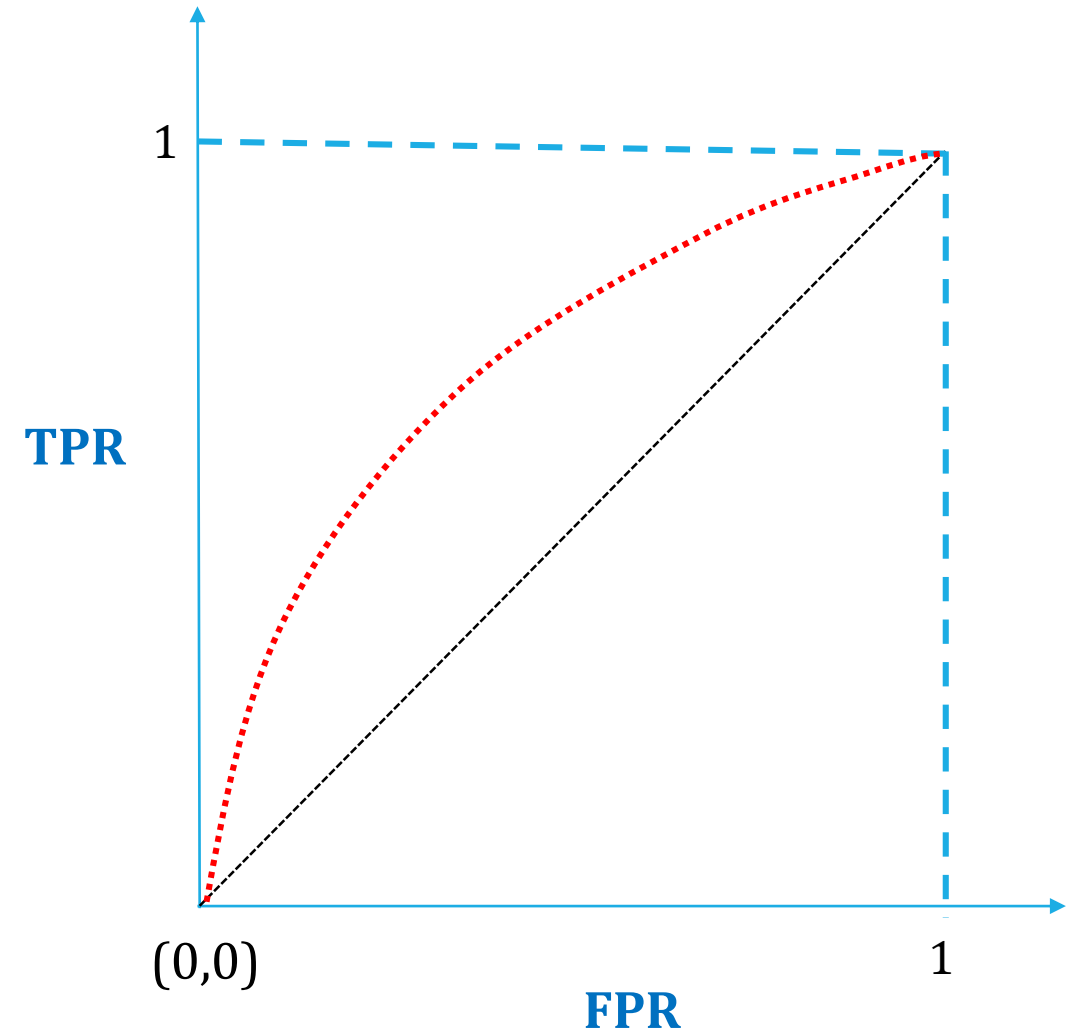
$$\text{FPR} = \text{False Positive Rate} = \frac{\text{FP}}{\text{FP} + \text{TN}} = \text{Probability of false alarm} = 1 - \text{specificity}$$

ROC Curve & AUC

ROC Curve = Receiving
Operating Characteristic Curve

ROC is a plot of FPR vs. TPR
calculated at different threshold
values.

Area **U**nder the ROC **C**urve is
called **AUC**



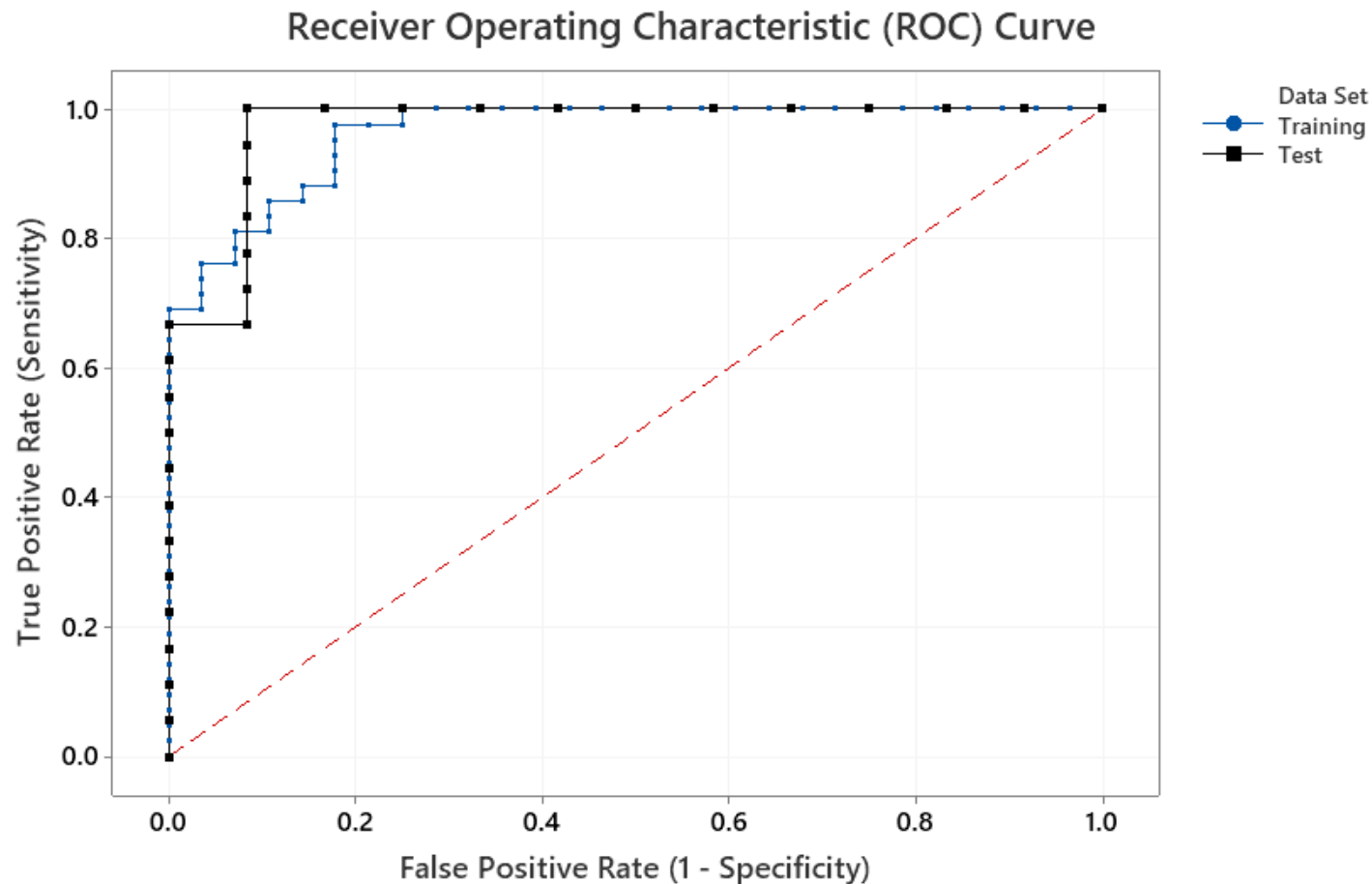
What is a good fit?

ROC is above the FPR = TPR line

AOC is more than 50%

ROC and AOC for Test data is close to that of Training data.

Large data set - Validation



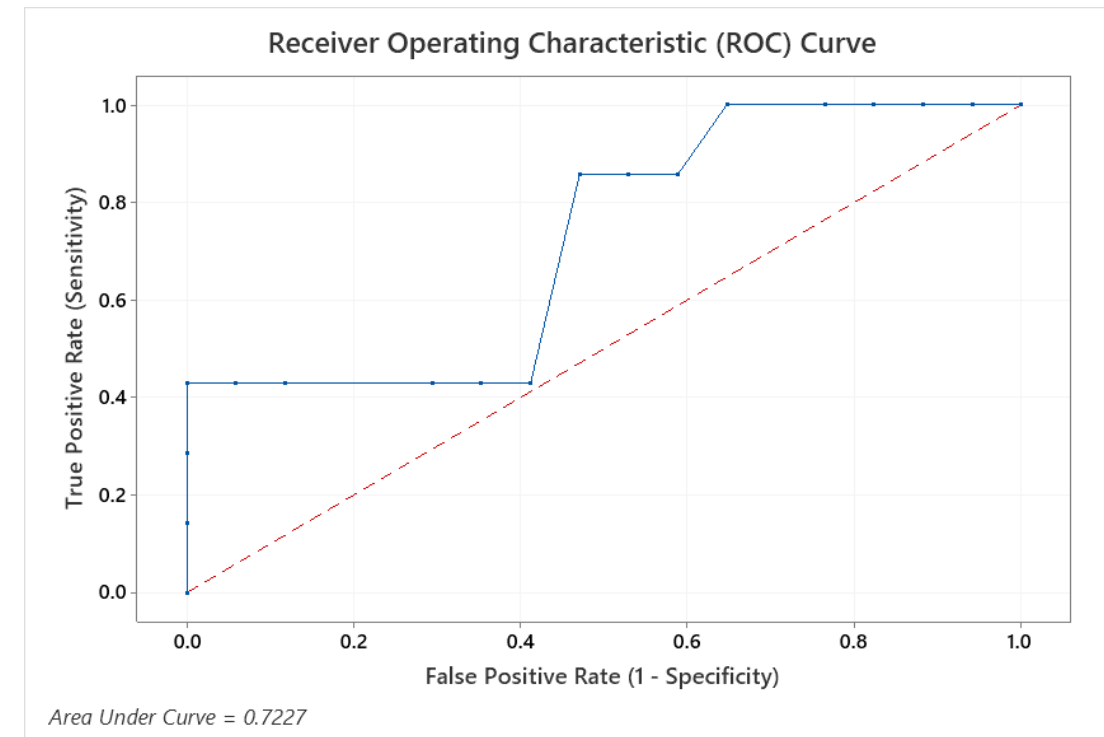
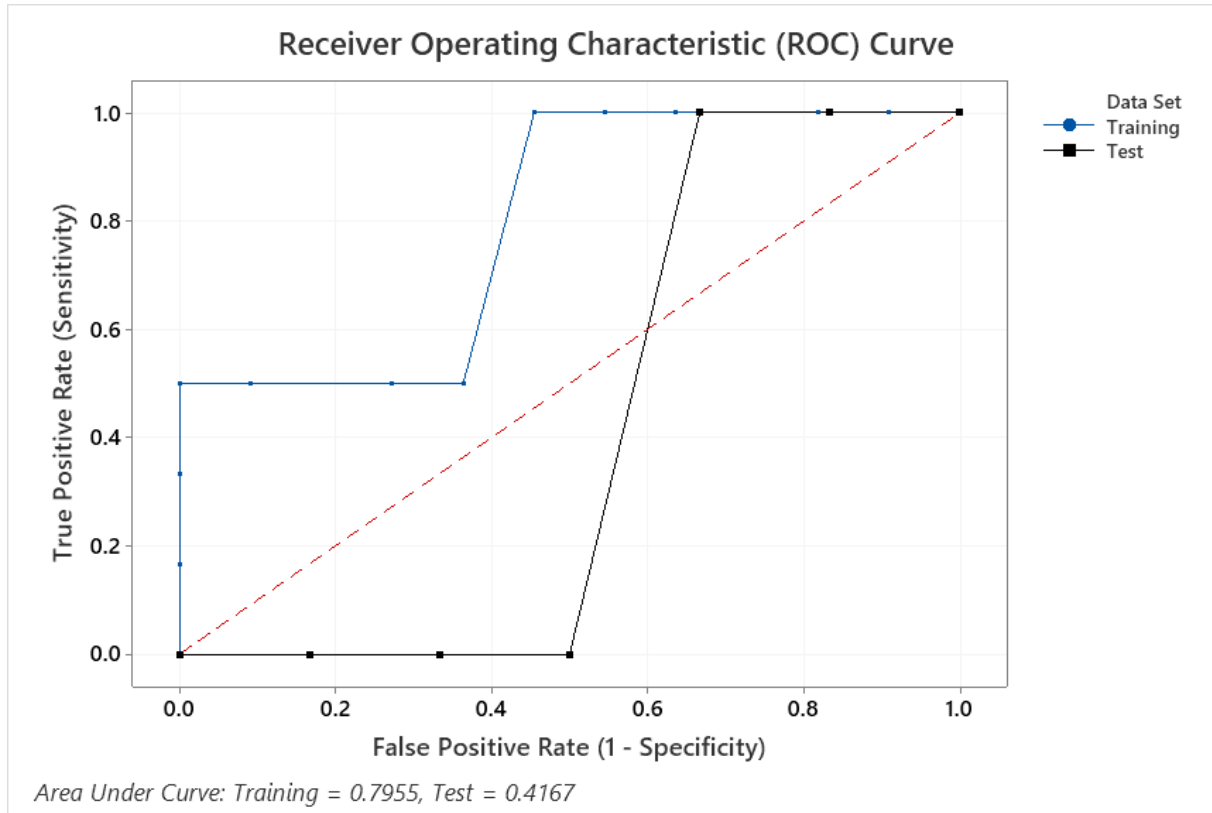
Validation

If ROC and AUC for Training Data and Test data are close then the model is validated.

If not valid:

- Too many features? Reduce the number of features.
- Correlated features? Keep only uncorrelated features.
- Training and test data sets chosen randomly? If not, correct it.
- Is data too small? Then avoid dividing the data. Keep only training data for estimation and make sure that the model fits well.

Small data – Caution!



Summary

Algorithm for Logistic Regression and choice of λ .

Validation of the model : Confusion Matrix

- ROC – Receiving Operating Characteristic Curve
- AUC – Area under the ROC Curve

Thank you...