

1 k-fache Kreuzvalidierung

Bei der k-fachen Kreuzvalidierung werden die Trainingsdaten k -mal in Testdaten und Validierungsdaten getrennt. Die Klassifizierer werden dann mit den Trainingsdaten trainiert und mit den Testdaten getestet und danach wieder angepasst. Das Testen der Klassifizierer wird mit diesen Daten vollzogen und danach wird das beste Modell ausgewählt.

2 Die Lernkurve

Bei der Lernkurve werden die Korrektklassifizierungsrate des Modells als y gegenüber der Anzahl der Daten. Bei der gewünschten Klassifizierungsrate als Konstante kann es sein, dass es einen sehr großen Bias gibt, das heißt die Korrektklassifizierungsrate der Test-, sowie der Validierungsdaten pendeln sich weit unter der Konstante der Klassifizierungsrate ein. Bei einer sehr hohen Varianz scheinen die Kurven beider sehr langsam gegen die angepeilte Rate zu konvergieren, das heißt das Modell läuft die Gefahr einer Überanpassung. Scheinen beide Kurven schnell gegen die angepeilte Korrektklassifizierungsrate zu konvergieren hat man einen guten Kompromiss zwischen Bias und Überanpassung gefunden.

3 Abstimmung durch Rastersuche

Bei Hyperparametern der Modelle kann durch reine Brute-Force-Methoden der optimale Hyperparameter für das Modell anhand der Korrektklassifizierungsrate bestimmt werden.

4 Wahrheitsmatrix

Bei einer Wahrheitsmatrix gibt es vier Fälle. Entweder es wurde so klassifiziert, wie es sein sollte oder es wurde das Gegenteil klassifiziert. Das ergibt vier Werte, die man sehr einfach auswerten kann. Die Korrektklassifizierungsrate berechnet sich dann konkret zu

$$\text{KKR} = \frac{\text{PR} + \text{PF}}{\text{NF} + \text{NR} + \text{PR} + \text{PF}} \quad (1)$$

Die Fehlerquote berechnet sich dann durch

$$\text{FQ} = 1 - \text{KKR} = \frac{\text{NF} + \text{NR}}{\text{NF} + \text{NR} + \text{PR} + \text{PF}} \quad (2)$$

Die Richtig-Positiv-Rate berechnet sich zu:

$$\text{RPR} = \frac{\text{PR}}{\text{NR} + \text{PR}} \quad (3)$$

Die Falsch-Positiv-Rate:

$$\text{FPR} = \frac{\text{PF}}{\text{PF} + \text{NF}} \quad (4)$$

Andere Werte zur Messung der Leistung sind die Genauigkeit

$$\text{GEN} = \frac{\text{PR}}{\text{NR} + \text{PR}} \quad (5)$$

Und die Trefferquote:

$$\text{TQ} = \frac{\text{PR}}{\text{NF} + \text{PR}} \quad (6)$$

Das $F1$ -Maß berechnet sich dann zu:

$$F1 = 2 \frac{\text{GEN} \times \text{TQ}}{\text{GEN} + \text{TQ}} \quad (7)$$

Für nichtbinäre Klassifizierer gelten dann die Formeln bei k Klassen:

$$\text{GEN}_{\text{mikro}} = \frac{\sum_{n=1}^k \text{PR}_n}{\sum_{n=1}^k (\text{NR}_n + \text{PR}_n)} \quad (8)$$

$$\text{GEN}_{\text{makro}} = \frac{\sum_{n=1}^k \text{Gen}_n}{k} \quad (9)$$