

Implementierung und Evaluation von Mixing-Modell-Alternativen für XCS

Markus Weis

1 EINLEITUNG

In modernen Michigan-Style Learning Classifier Systemen (LCS) wird üblicherweise aus einzelnen Classifier Vorhersagen eine Gesamtvorhersage berechnet. In XCS wird die Prediction über eine Fitness gewichtete Summe der einzelnen Predictions berechnet. Diese Fitness ist ein heuristisch Parameter, der inkrementell berechnet wird. Wilson schreibt, dass es wohl verschieden Möglichkeiten gibt, die Gesamtvorhersage zu berechnen, behält aber dieses Modell ohne weitere Ausführung in seinen späteren Arbeiten bei [4]. Drugowitsch setzt in seinem Buch "Design and Analysis of Learning Classifier Systems - A Probabilistic Approach" [1] dort an. Er beschreibt verschiedene Mixing Modelle und wertet diese empirisch aus. Einen dieser Modelle wird in der vorliegenden Arbeit beschrieben und eine Implementierung versucht: Inverse Variance Mixing. Aus den heuristisch Modellen schnitt dieses stets am besten ab, während die Laufzeit der analytischen Lösung nicht angemessen ist. Das Modell wird im Open-Source Projekt scikit-XCS [3], einer scikit-learn Erweiterung für XCS, implementiert.

2 INVERSE VARIANCE MIXING

Die Definition des Inverse Variance Mixing ist [1] entnommen.

$$\tau_k^{-1} = (c_k - D^X) \sum_n m_k(x_n) (\hat{w}_k^T x_n - y_n)^2$$

Dabei ist τ_k^{-1} die Varianz und folglich τ_k die inverse Invarianz. C_k ist der Matchcount des Classifiers k, D^X ist die Dimension des Input Raums. $m_k(x)$ ist eine Matching-Funktion für die gilt: $m_k(x) = 1$, falls die Bedingung des Classifiers k auf x passt und $m_k(x) = 0$ sonst. Der hintere Term stellt das quadratische Fehlermaß der Regression dar. Dieser muss hier angepasst werden, da Klassifikation implementiert werden soll.

Die erste Implementation lautet wie folgt:

$$\tau_k^{-1} = (c_k - D^X) \sum_n m_k(x_n) f_k(x_n)$$

Wobei $f_k(x)$ den Klassifikationsfehler darstellen soll: $f_k(x) = 0$, falls Classifier k die richtige Klasse für Input x hat und $f_k(x) = 1$, sonst.

Später wurde statt der Klasse noch der Fehler der Prediction als Fehlermaß verwendet:

$$\tau_k^{-1} = (c_k - D^X) \sum_n m_k(x_n) |p_k - P_k(x)|$$

p_k ist die Prediction des Classifiers k, $P_k(x)$ ist der Reward der ausgeschüttet wird wenn Classifier k auf Input x seine Klasse vorhersagt. Statt dem Betrag wurde auch versucht den Fehler zu quadrieren, dies hatte aber keine relevanten Auswirkungen.

Schließlich wird die inverse Varianz nicht direkt benutzt, sondern noch der sog. Gating Parameter g_k über das aktuelle Matchset berechnet:

$$g_k(x) = \frac{m_k(x) \tau_k}{m_i(x) \sum_i^K \tau_i}$$

Dadurch wird sichergestellt, dass die Summe aller Gating Parameter für jeden Input 1 ergibt.

3 PROBLEME

Im Buch wurde Inverse Variance Mixing nie im XCS-Context benutzt. Stattdessen wurde sie im allgemeineren Context von LCS definiert. LCS Classifier im Buch wurden durch Batch Learning trainiert werden und nicht iterativ, wie es bei XCS der Fall ist. Ein erster Ansatz der Implementierung hier war einfach die Fitness der einzelnen Classifier durch die inverse Variance zu ersetzen. Daraus ergaben sich aber 2 Probleme:

1. Da man für alle Classifier, die für die aktuelle Prediction verantwortlich sind, die Fitness updatet, müsste man für alle diese die inverse variance neu berechnen. Dazu wird aber für jeden Classifier das komplette Trainingsset durchlaufen. Dies führt zu einer nicht akzeptablen Laufzeit. Hier müsste wie im restlichen XCS ein iterativer Ansatz verfolgt werden. Dies wurde hier nicht gemacht, da man sich dadurch noch weiter von der ursprünglichen Idee entfernt und weiter Hyperparameter eingefügt werden müssten. Stattdessen wurde hier die Laufzeit erstmal in Kauf genommen und geschaut, ob auf diese Weise gute Ergebnisse auftreten.

2. Das Erzeugen neuer Classifier benutzt die Fitness als Auswahlkriterium. Allerdings ist die Inverse Variance zu Beginn des Training Vorgangs noch nicht als sinnvoll zu betrachten. Beispielsweise wird ein Classifier, der auf einen Input perfekt passt stets bevorzugt, egal ob es allgemeinere gibt, die ähnlich gut sind.

Wir sind schließlich zu dem Schluss gekommen, die Fitness des ursprünglichen XCS Systems in Ruhe zu lassen und die Inverse Variance nur für das Mixing zu benutzen. Das Problem mit der Laufzeit besteht aber, da für die Prediction bei nicht Explorationen berechnet werden muss. Schließlich wurden 2 Alternativen implementiert. Die eine macht die Prediction während des Trainings bereits mit der inversen Invarianz und zeigt deshalb eine sehr schlechte Laufzeit auf. Bei der anderen wird die neue Mixing Art nur bei Predictions benutzt, die nicht während des Trainings auftreten, also nur im fertigen Modell. Das hat den Vorteil, dass die inverse Varianz für alle Classifier am Ende des Trainingsprozesses berechnet werden kann und für die einzelnen Prediction nur der Gating Parameter aktualisiert werden muss.

Wie bereits beschrieben handelt es sich bei der eigentlichen Definition um Regression. Bei Regression ist es unwahrscheinlich bzw. unmöglich, je nach Modell, auf einen perfekten Classifier zu treffen. Das Problem bei einem perfekten Classifier ist, dass der hintere Teil der Formel, also das Fehlermaß, Null ergibt. Das wird zum

	6 Bit	11 Bit	20 Bit
normal	0,98 (+/- 0,03)	0,98 (+/- 0,02)	0,89 (+/- 0,05)
evenly-distributed	0,88 (+/- 0,06)	0,92 (+/- 0,04)	0,84 (+/- 0,05)
only-mixing	0,87 (+/- 0,06)	0,90 (+/- 0,04)	0,84 (+/- 0,05)
continuous-update	0,93 (+/- 0,04)	0,93 (+/- 0,03)	0,85 (+/- 0,05)

Tabelle 1: Durchschnittliche Accuracy der verschiedenen Implementierung mit Standartabweichung in Klammern.

Problem, wenn man zur Berechnung der Inversen Varianz den Umkehrbruch bildet. Im Buch ist das nicht weiter definiert. In dieser Implementation habe ich in diesem Fall Unendlich (numpy.inf) eingesetzt. Dadurch bekommt man schließlich im Mixing Probleme. Zwar bekommt man bei einem Unendlich beim Mixing-Parameter g_k bei allen anderen Classifiern, die nicht Unendlich sind 0 raus, allerdings ist $\frac{\text{numpy.inf}}{\text{numpy.inf}}$ nicht definiert. Stattdessen wird g_k im Falle eines Classifiers k im Matchset mit $\gamma_k = \infty$ wie folgt berechnet: $g_k = \frac{1}{c_{inf}}$ für alle Classifier k mit $\gamma_k = \infty$ und $g_k = 0$ für alle anderen Classifier. c_{inf} ist die Anzahl der Classifier mit $\gamma_k = \infty$.

Ein weiteres Problem ergab sich im Fall, dass der Matchcount c_k für einige Classifier geringer als die Dimension des Inputs war, sodass sich für die Varianz und die negative Varianz ein negativer Wert ergibt. Im Buch ist nicht beschrieben, wie in diesem Fall mit den Classifiern umgegangen werden soll. In meiner Implementation werden diese Classifier einfach ignoriert.

4 ERGEBNISSE

Die Implementierung wurde auf Multiplexern verschiedener Größe evaluiert: 6 Bit, 11 Bit und 20 Bit. Die Trainingsdauer und die Hyperparameter orientieren sich an [2]. Dort ist zu entnehmen, wie viele Iterationen XCS braucht, um die verschiedenen großen Multiplexer perfekt zu lernen. Letztendlich wäre jedes der hier zu evaluierenden Implementierungen bei den Problemen zu einer perfekten Accuracy gekommen, wenn man nur lange genug trainiert. Es wäre wohl sinnvoll die Implementierungen auch auf Problemen zu testen, bei denen XCS keine perfekte Lösung finden kann. Dies wird hier allerdings nicht gemacht. Es wird für 4000, 10000, 40000 Iterationen für den 6-Bit, 11-Bit und 20-Bit Multiplexer trainiert. Die durchschnittliche Accuracy der Cross-Validation der verschiedenen Implementierungen ist der Tabelle ?? zu entnehmen. Dabei wurde beim 6-Bit Multiplexer eine 10x5, beim 11-Bit Multiplexer eine 5x5 und beim 20-Bit Multiplexer eine 2x5 Cross-Validation durchgeführt. Dies liegt an der langen Laufzeit der Implementierung mit dem kontinuierlichen Update der inversen Varianz. Es wird neben den beiden bereits erwähnten Implementierungen der inversen Varianz noch eine Gleichgewichtung der Klassifier vorgenommen, um den Effekt des Mixings vergleichen zu können (evenly-distributed).

Es ist hier bereits zu erkennen, dass die Ergebnisse des inverse Variance Mixing schlechter sind als die des normalen XCS-Mixing. Genauere Auswertungen der verschiedenen Problem-Größen sind hier aufgeführt.

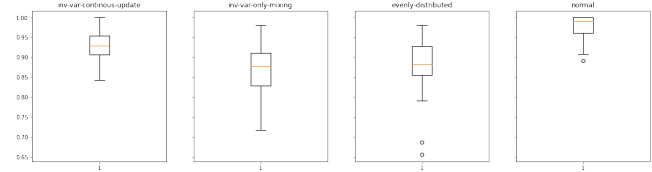


Abbildung 1: Accuracy Werte des 6 Bit Multiplexers

	inv-var-continuous-update	inv-var-only-mixing	evenly-distributed	normal
inv-var-continuous-update	0.000000	0.004355	0.034659	0.000425
inv-var-only-mixing	0.000000	0.000000	0.559876	0.047220
evenly-distributed	0.000000	0.000000	0.000000	0.001125
normal	0.000000	0.000000	0.000000	0.000000

Abbildung 2: 5x2 Cross-validation paired ttest auf dem 6 Bit Multiplexer

4.1 6 Bit Multiplexer

Die Boxplots zu den Accuracy Werten aus Tabelle ?? sind in Abbildung 5 zu sehen. Es ist deutlich, dass das normale XCS am Besten abschneidet. Evenly-distributed ist zwar auf ähnlichem Niveau, wie die beiden Inverse Variance Arten, allerdings hat es einige negative Ausreißer.

Schließlich wurde noch ein TTest durchgeführt, um zu überprüfen, ob die Unterschiede der Algorithmen statische Signifikanz aufweisen. Hierfür wurde die Funktion `paired_ttest_5x2cv()` der MLxtend Library von Sebastian Raschka verwendet. Es wurde ein Signifikanzwert von $\alpha = 0,05$ verwendet. Dass heißt die Unterschiede zwischen den Implementierungen sind signifikant, fass das Ergebnis der des `paired_ttest_5x2cv()` einen Wert unter 0,05 liefert. Diese Einträge sind gelb hinterlegt. Die Ergebnisse sind Abbildung 4 zu entnehmen.

Der TTest erkennt Unterschiede zwischen allen Verteilungen, außer zwischen evenly-distributed und inv-var-only-mixing. Ein Ergebnis, das wieder gegen die Verwendung dieser Implementierung spricht.

4.2 11 Bit Multiplexer

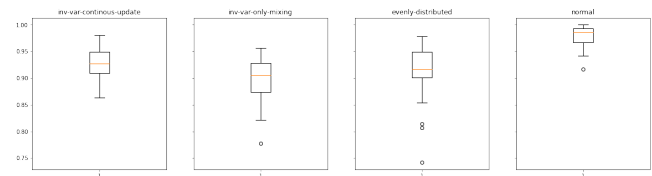


Abbildung 3: Accuracy Werte des 6 Bit Multiplexers

4.3 20 Bit Multiplexer

5 FAZIT

Die hier beschriebenen Ergebnisse zeigen, dass das Inverse Variance Mixing, wie es hier implementiert ist, sicher nicht als sinnvolle Alternative zum gewöhnlichen XCS-Mixing benutzt werden kann. In

	inv-var-continous-update	inv-var-only-mixing	evenly-distributed	normal
inv-var-continous-update	0.000000	0.002321	0.014073	0.000045
inv-var-only-mixing	0.000000	0.000000	0.011643	0.000222
evenly-distributed	0.000000	0.000000	0.000000	0.000037
normal	0.000000	0.000000	0.000000	0.000000

Abbildung 4: 5x2 Cross-validation paired ttest auf dem 6 Bit Multiplexer

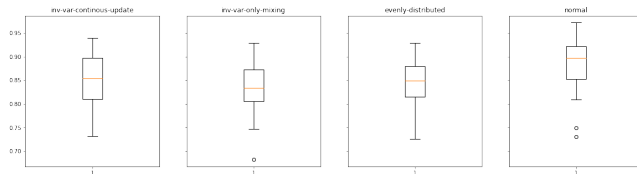


Abbildung 5: Accuracy Werte des 6 Bit Multiplexers

dieser Implementierung wurde allerdings notwendigerweise stark von der eigentlichen vorgeschlagenen Heuristik abgewichen. Welche der Abweichungen zu den schlechten Resultaten führt ist schwer

zu sagen. Es handelte sich hierbei eher um einen praktischen Versuch und die formale Korrektheit wurde dabei vernachlässigt. Dies war aufgrund der Komplexität des Themas und der beschränkten Zeit notwendig. Es kann also sein, dass es durchaus sinnvolle Implementierung von Inverse Variance Mixing in XCS gibt, v.a. da hier nur Klassifikation betrachtet wurde. Auffällig ist, dass das Mixing mit gleicher Verteilung der Classifier relativ gut abschneidet und bei hier vorliegenden Versuchsaufbau keine statische Signifikanz zum normalen XCS-Mixing vorliegt. Dies deutet darauf hin, dass das XCS Mixing nicht optimal ist. Weiter Mixing Modelle sollten daher ausprobiert werden.

LITERATUR

- [1] Jan Drugowitsch. 2008. *Design and Analysis of Learning Classifier Systems - A Probabilistic Approach*. Vol. 139. <https://doi.org/10.1007/978-3-540-79866-8>
- [2] Muhammad Iqbal, Will Browne, and Mengjie Zhang. 2013. Learning complex, overlapping and niche imbalance Boolean problems using XCS-based classifier systems. *Evolutionary Intelligence* 6 (11 2013). <https://doi.org/10.1007/s12065-013-0091-1>
- [3] UrblsLab. 2020. scikit-XCS. *Github* (2020). <https://github.com/UrblsLab/scikit-XCS>
- [4] Stewart W. Wilson. 1995. Classifier Fitness Based on Accuracy. *Evol. Comput.* 3, 2 (June 1995), 149–175. <https://doi.org/10.1162/evco.1995.3.2.149>