

3 Bias-Variance und Regularisierung

Systemtechnik BSc
FS 2025

Aufgaben

Applied Neural Networks im Modul ANN | WUCH

Lernziele

Nach dem Bearbeiten dieser Übungsserie ...

- können Sie mit `pytorch-lightning` einfache sequentielle neuronale Netze für die Klassifizierung und Regression aufbauen.

Aufgaben:

- können Sie beurteilen, ob ein NN-Modell eine hohe Varianz (overfitting) oder eher einen hohen Bias (underfitting) aufweist. Sie sind in der Lage, mit geeigneten Regularisierungsmassnahmen eine Überanpassung eines Modells zu vermeiden (weight regularization [L1,L2], early stopping, dropout, model complexity)

Aufgaben: 1

Über- und Unteranpassung

Overfitting

Overfitting bedeutet: Unser Modell passt sich den Trainingsdaten *zu sehr* an. Es ist wichtig zu lernen, wie man mit der Überanpassung umgeht. Obwohl es oft möglich ist, eine hohe Genauigkeit im *Trainingsdatensatz* zu erreichen, wollen wir eigentlich Modelle entwickeln, die sich gut auf einen *Testdatensatz* verallgemeinern lassen (oder auf Daten, die sie noch nie gesehen haben).

Underfitting

Das Gegenteil von Overfitting ist *Underfitting*. Underfitting liegt vor, wenn die Performance des Modells auf den Trainingsdaten noch verbesserungswürdig sind. Dies kann aus verschiedenen Gründen geschehen: Wenn das Modell nicht leistungsfähig genug ist, überreguliert ist oder einfach nicht lange genug trainiert wurde. Das bedeutet, dass das Netz die relevanten Muster in den Trainingsdaten nicht gelernt hat.

Wenn Sie jedoch zu lange trainieren, wird das Modell anfangen, sich zu sehr anzupassen und Muster aus den Trainingsdaten zu lernen, die sich nicht auf die Testdaten übertragen lassen. Wir müssen also ein Gleichgewicht finden. Es ist nützlich zu wissen, wie man für eine angemessene Anzahl von Epochen trainiert, wie wir weiter unten erläutern werden.

Regularisierung

Um eine Überanpassung zu vermeiden, ist die beste Lösung die Verwendung vollständigerer Trainingsdaten. Der Datensatz sollte die gesamte Bandbreite der Eingaben abdecken, die das Modell verarbeiten soll. Zusätzliche Daten sind nur dann sinnvoll, wenn sie neue und interessante Fälle abdecken. Der einfachste Weg, eine Überanpassung zu verhindern, ist, mit einem kleinen Modell zu beginnen: Ein Modell mit einer kleinen Anzahl von lernbaren Parametern (die durch die Anzahl der Schichten und die Anzahl der Einheiten pro Schicht bestimmt wird). Beim Deep Learning wird die Anzahl der lernbaren Parameter in einem Modell oft als *Kapazität* des Modells bezeichnet. Deep-Learning-Modelle neigen dazu, sich gut an die Trainingsdaten anzupassen, aber die eigentliche Herausforderung ist die Verallgemeinerung, nicht die Anpassung.

Wenn das Netzwerk andererseits nur über begrenzte Speicherressourcen verfügt, kann es das Mapping nicht so leicht erlernen. Um seinen Verlust zu minimieren, muss es komprimierte Darstellungen lernen, die eine höhere Vorhersagekraft haben. Wenn Sie Ihr Modell jedoch zu klein machen, wird es Schwierigkeiten haben, sich an die Trainingsdaten anzupassen. Es gibt ein Gleichgewicht zwischen „zu viel Kapazität“ und „nicht genug Kapazität“.

Ein Modell, das auf umfangreicheren Daten trainiert wurde, wird natürlich besser verallgemeinern. Wenn dies nicht mehr möglich ist, besteht die nächstbeste Lösung darin, Techniken wie die *Regularisierung* anzuwenden. Diese schränken die Menge und Art

der Informationen ein, die Ihr Modell speichern kann. Wenn sich ein Netzwerk nur eine kleine Anzahl von Mustern merken kann, wird es durch den Optimierungsprozess gezwungen, sich auf die auffälligsten Muster zu konzentrieren, die eine bessere Chance haben, gut zu verallgemeinern. In dieser Übungsaufgabe werden wir verschiedene *gängige Regularisierungstechniken* untersuchen und sie zur Verbesserung eines Klassifizierungsmodells einsetzen.

Aufgabe 1. Benchmark: Über- und Unteranpassung

- (a) Laden Sie den Datensatz `HIGGS.csv.gz` von `mlphysics.ics.uci.edu` herunter. Achtung: Es sind rund 2.7 GB an Daten. Legen Sie den Datensatz auf einen Ordner Ihrer Wahl ab. Sie brauchen das komprimierte File nicht zu entzippen.

Die Daten wurden mit Hilfe von Monte-Carlo-Simulationen erstellt. Die ersten 21 Merkmale (Spalten 2-22) sind kinematische Eigenschaften, die von den Teilchendetektoren im Beschleuniger gemessen wurden. Die letzten sieben Merkmale sind Funktionen der ersten 21 Merkmalen; es handelt sich dabei um hochrangige Merkmale, die von Physikern abgeleitet wurden, um zwischen den beiden Klassen zu unterscheiden. Das Ziel dieser Übungsserie ist nicht die Teilchenphysik, daher sollten Sie sich nicht mit den Details des Datensatzes beschäftigen. Er enthält 11'000'000 Beispiele, jedes mit 28 Merkmalen und einem binären Klassenlabel.

- (b) Öffnen Sie das Jupyter-Template `overfit_and_underfit_TEMPLATE.ipynb` und führen Sie die erste Code-Zelle aus, um die Daten zu laden. Passen Sie den Datenpfad entsprechend an. Beim Setup kann das Package `ipywidgets` installiert werden für schönere Ladebalken während des Trainings. Es enthält folgende Funktionen:

download_file_if_not_exists Laden der Daten

config_dict für allgemeine Einstellungen

HiggsDataset Klasse für die Daten

create_model für das einfache Erstellen der Architektur

LightningModel Klasse für das Training

train_model allgemeine Einstellungen für das Training für mehr Übersichtlichkeit

train_with_early_stopping Early Stopping, um abzubrechen, falls das Modell nicht mehr verbessert

- (c) **Winziges Modell:** Um eine geeignete Modellgrösse zu finden, beginnen Sie am besten mit relativ wenigen Schichten und Parametern und beginnen dann, die Grösse der Schichten zu erhöhen oder neue Schichten hinzuzufügen, bis Sie eine Verringerung des Validierungsverlustes feststellen. Beginnen Sie mit einem einfachen Modell, das nur nn.Linear-Layers als Basis verwendet, erstellen Sie dann grössere Versionen und vergleichen Sie diese. Im Code ist eine config_dict in der die Layers definiert werden.

```

1 config_tiny = {
2     "model": {
3         "input_size": 28, # Eingangsgroesse des Modells
4         "tiny": {
5             "hidden_layers": [16],
6             "dropout_rate": 0.0,
7             "l2_reg": 0.0,
8         }, # Konfiguration fuer Tiny-Modell
9     },
10 }
11
```

Trainieren Sie das kleine Modell mit folgenden Schritten und es wird auch direkt der Trainingsloss und Validierungsloss geplottet.

```

1 print("Training Basis Tiny Model...")
2 # printe die Konfiguration des Modells
3 print(config_tiny["model"]["tiny"])
4 tiny_model = create_model(
5     config_tiny["model"]["input_size"], **config_tiny["model"]["tiny"]
6 )
7 tiny_lightning_model = LightningModel(tiny_model)
8 tiny_trained = train_with_early_stopping(
9     tiny_lightning_model, train_loader, validate_loader
10 )
11
12
```

- (d) **Kleines Modell:** Um zu sehen, ob Sie die Leistung des kleinen Modells übertreffen können, trainieren Sie nach und nach einige grössere Modelle. Versuchen Sie es mit zwei versteckten Schichten mit je 16 Einheiten:

```
1 config_small = {
2     "model": {
3         "input_size": 28, # Eingangsgroesse des Modells
4         "small": {
5             "hidden_layers": [16, 16],
6             "dropout_rate": 0.0,
7             "l2_reg": 0.0,
8         }, # Konfiguration fuer Small-Modell
9     },
10 }
11
```

Trainiert wird das Modell wie folgt:

```
1 print("Training Small Model (Basis, [16,16])...")
2 # printe die Konfiguration des Modells
3 print(config_small["model"]["small"])
4 small_model = create_model(
5     config_small["model"]["input_size"], **config_small["model"]["small"]
6 )
7 small_lightning_model = LightningModel(small_model)
8 small_trained = train_with_early_stopping(
9     small_lightning_model, train_loader, validate_loader
10 )
11
```

- (e) **Mittleres Modell:** Testen Sie nun ein Modell mit drei versteckten Schichten (hidden layers) mit je 64 Einheiten.

```
1 config_medium = {
2     "model": {
3         "input_size": 28, # Eingangsgroesse des Modells
4         "medium": {
5             "hidden_layers": [64, 64, 64],
6             "dropout_rate": 0.0,
7             "l2_reg": 0.0,
8         }, # Konfiguration fuer Medium-Modell
9     },
10 }
11
```

Auch dieses Modell wird wie folgt Trainiert.

```
1 print("Training Medium Model (Basis, [64,64,64])...")
2 # printe die Konfiguration des Modells
3 print(config_medium["model"]["medium"])
```

```

4 medium_model = create_model(
5 config_medium["model"]["input_size"], **config_medium["model"]["medium"]
6 )
7 medium_lightning_model = LightningModel(medium_model)
8 medium_trained = train_with_early_stopping(
9 medium_lightning_model, train_loader, validate_loader
10 )
11

```

- (f) **Grosses Modell:** Als Nächstes fügen wir zu diesem Benchmark ein Netzwerk hinzu, das viel mehr Kapazität hat, weit mehr als das Problem rechtfertigen würde:

```

1 config_large = {
2     "model": {
3         "input_size": 28, # Eingangsgroesse des Modells
4         "large": {
5             "hidden_layers": [512, 512, 512, 512],
6             "dropout_rate": 0.0,
7             "l2_reg": 0.0,
8         }, # Konfiguration fuer Large-Modell
9     },
10 }
11
12 print("Training Large Model (Basis, [512,512,512,512])...")
13 # printe die Konfiguration des Modells
14 print(config_large["model"]["large"])
15 large_model = create_model(
16 config_large["model"]["input_size"], **config_large["model"]["large"]
17 )
18 large_lightning_model = LightningModel(large_model)
19 large_trained = train_with_early_stopping(
20 large_lightning_model, train_loader, validate_loader
21 )
22

```

- (g) **Interpretation:** Es werden alle Lernkurven zum Modell auch geplottet, beschreiben Sie in wenigen Sätzen, ob eine Über oder Unteranpassung erkennbar ist Anhand des Trainings- und Validierungslosses.
- (h) **Weight Regularization:** Vielleicht kennen Sie das Prinzip von OCCAMS Rasiermesser: Wenn es zwei Erklärungen für etwas gibt, ist die Erklärung, die am wahrscheinlichsten richtig ist, die „einfachste“, diejenige, die die wenigsten Annahmen

enthält. Dies gilt auch für die Modelle, die von neuronalen Netzen gelernt werden: Bei bestimmten Trainingsdaten und einer Netzarchitektur gibt es mehrere Sätze von Gewichtungswerten (mehrere Modelle), die die Daten erklären könnten, und bei einfacheren Modellen ist die Wahrscheinlichkeit einer Überanpassung geringer als bei komplexen Modellen.

Ein „einfaches Modell“ ist in diesem Zusammenhang ein Modell, bei dem die Verteilung der Parameterwerte eine geringere *Entropie* aufweist (oder ein Modell mit insgesamt weniger Parametern, wie wir im obigen Abschnitt gesehen haben). Eine gängige Methode zur Abschwächung der Überanpassung besteht daher darin, die Komplexität eines Netzes einzuschränken, indem seine Gewichte gezwungen werden, nur kleine Werte anzunehmen, wodurch die Verteilung der Gewichtswerte „regelmässiger“ wird. Dies wird als *Gewichtsregularisierung* bezeichnet und erfolgt durch Hinzufügen von Kosten zur Verlustfunktion des Netzes, die mit grossen Gewichten verbunden sind. Diese Kosten gibt es in zwei Varianten:

- **L1-Regularisierung**, bei der die hinzugefügten Kosten proportional zum absoluten Wert der Gewichtskoeffizienten sind (d. h. zur so genannten L1-Norm der Gewichte).
- **L2-Regularisierung** bei der die zusätzlichen Kosten proportional zum Quadrat des Wertes der Gewichtungskoeffizienten sind (d.h. zu dem, was man die quadrierte L2-Norm der Gewichte nennt). Die L2-Regularisierung wird im Zusammenhang mit neuronalen Netzen auch als Gewichtsabnahme bezeichnet.

Die L1-Regularisierung verschiebt einige Gewichte genau gegen Null, was ein spärliches Modell fördert. Die L2-Regularisierung bestraft die Parameter der Gewichte, ohne sie spärlich zu machen, da die Strafe bei kleinen Gewichten gegen Null geht – ein Grund, warum L2 häufiger verwendet wird.

In PyTorch Lightning wird die L2-Gewichtsregulierung (auch bekannt als Weight Decay) typischerweise über den Optimierer hinzugefügt.

```
1 def configure_optimizers(self):  
2     optimizer = optim.Adam(self.parameters(), lr=self.lr, weight_decay=  
    self.l2_reg) # Optimierer  
3
```

- (i) **Dropout:** Dropout ist eine der effektivsten und am häufigsten verwendeten Regularisierungstechniken für neuronale Netze, die von Hinton und seinen Studenten an der Universität von Toronto entwickelt wurde. Die intuitive Erklärung für Dropout ist, dass sich einzelne Knoten im Netz nicht auf die Ausgaben der anderen verlassen können, sondern dass jeder Knoten für sich selbst nützliche Merkmale ausgeben muss. Dropout, angewandt auf eine Schicht, besteht darin, dass während des Trainings eine Anzahl von Ausgangsmerkmalen der Schicht zufällig „weggelassen“ (d. h. auf Null gesetzt) wird. Nehmen wir an, eine gegebene Schicht hätte normalerweise einen Vektor $[0.2, 0.5, 1.3, 0.8, 1.1]$ für eine gegebene Eingabeprobe während des Trainings geliefert; nach Anwendung von Dropout wird dieser Vektor einige zufällig verteilte Nulleinträge haben, z. B. $[0, 0.5, 1.3, 0, 1.1]$.

Die *Dropout-Rate* ist der Anteil der Merkmale, die mit Nullen versehen werden; sie wird normalerweise zwischen 0.2 und 0.5 festgelegt. Zur Testzeit werden keine Einheiten herausgenommen, stattdessen werden die Ausgabewerte der Schicht um den Faktor der Dropout-Rate herunterskaliert, um die Tatsache auszugleichen, dass mehr Einheiten aktiv sind als zur Trainingszeit.

Über das `config_` Dictionary kann die `dropout_rate` mitgegeben werden. Sie wird dann bei `create_model` hinzugefügt und wie in der Funktion `create_model` definiert erstellt.

```
1 def create_model(input_size, hidden_layers, dropout_rate=0.0, l2_reg
  =0.0):
2     layers = []
3     prev_dim = input_size
4     for layer_dim in hidden_layers:
5         layers.append(nn.Linear(prev_dim, layer_dim))
6         layers.append(nn.ELU())
7         if dropout_rate > 0:
8             layers.append(nn.Dropout(dropout_rate))
9         prev_dim = layer_dim
10    layers.append(nn.Linear(prev_dim, 1))
11    return nn.Sequential(*layers)
12
```

- (j) **Kombinierte Anwendung von L2 und Dropout:** In diesem Schritt werden L2 und Dropout noch zusammen kombiniert.

- (k) **Spielen Sie mit verschiedenen Parametern und Kombinationen:** Bewerten Sie die Kurven. Versuchen Sie das beste Modell zu finden, welche Hyperparameter hat dieses Modell?

Lösungen

Lösung 1.

overfit_and_underfit_SOLUTION-PyTorch.ipynb