



1 Linear Least Squares Regression (12 Punkte)

In der Vorlesung haben wir zwei Arten von Modellen zur Systemidentifikation gesehen: White-Box und Black-Box. Ein Typ von Black-Box-Modellen ist die *Linear Least Squares Regression*. In dieser Übung beschäftigen wir uns mit den Grundlagen der linearen Regression.

- a) Was sind die Unterschiede zwischen White-Box- und Black-Box-Modellen? Was sind die Vor- und Nachteile des einen oder des anderen?

Lösungsvorschlag: (1.0 Punkt(e)) **White-Box Modellen:** Nutzen unser physikalisches Wissen, um parametrische Funktion festzulegen.

- Vorteil: Wir können **Vorkenntnisse** in unser Modell **einbeziehen**; In einigen Fällen kann dies **die Komplexität des Erlernens der Parameter reduzieren**
- Nachteil: Was wir nicht modellieren, kann nicht repräsentiert werden; Un- oder falsch modellierte Aspekte können systematisch die Parameterschätzung verfälschen

Black-Box Modellen: Anstatt komplizierten Physik Effekte zu modellieren, kann man direkt mit Modellen $y = f(x; \theta)$ arbeiten, welches sehr viele verschiedene Funktionen darstellen kann

- Vorteil: Das **Training** dieser Modelle ist manchmal **einfacher** als White-Box-Modelle; **robuster gegenüber Rauschen bei den Messungen**
- Nachteil: Keine physikalischen Prinzipien werden verwendet um das Modell herzuleiten

- b) Ausgehend von einem parametrisierten linearen Modell $\hat{y} = \theta^T \psi(x)$ wollen wir die Parameter θ so finden, dass unsere Schätzung nahe an den Daten $D = \{(x_i, y_i)\}$ liegt, wo $x_i \in \mathbb{R}^{d_x}$, $y_i \in \mathbb{R}$, $\theta \in \mathbb{R}^{d_\theta}$ und $\psi: \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_\theta}$. Die optimalen Parameter θ^* können durch Minimierung der Funktion der *Least Squares Error* gefunden werden

$$\theta^* = \arg \min_{\theta} L(\theta) = \arg \min_{\theta} \frac{1}{N} \left(y - \psi^T \theta \right)^T \left(y - \psi^T \theta \right), \quad (1)$$

mit den gesammelten Zieldaten $y \in \mathbb{R}^N$, und der *Feature Matrix* $\psi \in \mathbb{R}^{d_\theta \times N}$, die in den Spalten die Ausgabe der *Feature function* $\psi(x)$ für jeden Datenpunkt enthält.

Finde analytisch den Parametervektor θ , der die *Least Squares Error* $L(\theta)$ in (1) minimiert.

Hinweise: $(AB)^T = B^T A^T$. $\nabla_z (z^T A z) = 2Az$, falls A quadratisch und symmetrisch ist. $\nabla_z (z^T a) = \nabla_z (a^T z) = a$

Lösungsvorschlag: (2.0 Punkt(e))

$$\begin{aligned}L(\boldsymbol{\theta}) &= \frac{1}{N} (\mathbf{y} - \boldsymbol{\psi}^T \boldsymbol{\theta})^T (\mathbf{y} - \boldsymbol{\psi}^T \boldsymbol{\theta}) \\&= \frac{1}{N} (\mathbf{y}^T \mathbf{y} - \mathbf{y}^T \boldsymbol{\psi}^T \boldsymbol{\theta} - \boldsymbol{\theta}^T \boldsymbol{\psi} \mathbf{y} + \boldsymbol{\theta}^T \boldsymbol{\psi} \boldsymbol{\psi}^T \boldsymbol{\theta}) \\ \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) &= -2\boldsymbol{\psi} \mathbf{y} + 2\boldsymbol{\psi} \boldsymbol{\psi}^T \boldsymbol{\theta} = 0 \\ \boldsymbol{\theta}^* &= (\boldsymbol{\psi} \boldsymbol{\psi}^T)^{-1} \boldsymbol{\psi} \mathbf{y}\end{aligned}$$

- c) Wenn wir bei der Anpassung eines parametrisierten Modells an die Daten nicht vorsichtig sind, kann es passieren, dass das Modell sich zu sehr an die Daten anpasst (*overfit* auf Englisch), d.h. der durchschnittliche Fehler (*Mean Squared Error*) ist für die Trainingsdaten fast null, aber wenn du mit dem Modell eine Vorhersage in einem neuem Punkt \mathbf{x} machst, wird dein Ergebnis völlig vom wahren Ergebnis abweichen. Eine Möglichkeit, das *overfit* zu reduzieren, besteht darin, der Kostenfunktion $L(\boldsymbol{\theta})$ eine Strafe in den Parametern hinzuzufügen. Wir fügen eine Strafe proportional zur Norm von $\boldsymbol{\theta}$ hinzu (heißt *Ridge Regression*), $\lambda \boldsymbol{\theta}^T \boldsymbol{\theta}$ für $\lambda > 0$, die Kostenfunktion ist dann $L(\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}^T \boldsymbol{\theta}$. Finde analytisch mit dieser modifizierten Kostenfunktion wieder die neuen optimalen Parametern.

Lösungsvorschlag: (1.0 Punkt(e))

$$\begin{aligned}L(\boldsymbol{\theta}) &= \frac{1}{N} (\mathbf{y} - \boldsymbol{\psi}^T \boldsymbol{\theta})^T (\mathbf{y} - \boldsymbol{\psi}^T \boldsymbol{\theta}) + \lambda \boldsymbol{\theta}^T \boldsymbol{\theta} \\ \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) &= -\frac{2}{N} \boldsymbol{\psi} \mathbf{y} + 2\boldsymbol{\psi} \boldsymbol{\psi}^T \boldsymbol{\theta} + 2\lambda \boldsymbol{\theta} = 0 \\ \boldsymbol{\theta}^* &= (\boldsymbol{\psi} \boldsymbol{\psi}^T + \lambda N \mathbf{I})^{-1} \boldsymbol{\psi} \mathbf{y}\end{aligned}$$

- d) Bitte lade die Datei uebung09-programming.zip aus dem Moodle herunter.

Wir werden einen Roboter in einem unwegsamen Gelände einsetzen und möchten wissen, welche Kraft auf die Geschwindigkeit des Roboters wirkt. Wir gehen davon aus, dass es einen linearen Zusammenhang zwischen der resultierenden Kraft und der Geschwindigkeit des Roboters gibt (die Kraft kann ungleich Null sein, wenn sich der Roboter nicht bewegt, z.B. unter der Annahme, dass eine konstante äußere Kraft vorhanden ist). Wir bilden diese lineare Beziehung als Black-Box-Modell $\hat{y} = \boldsymbol{\theta}^T \boldsymbol{\psi}(\mathbf{x})$ ab, bei dem \mathbf{x} die Geschwindigkeit und y die Kraft sind. Wir sammeln N Messungen aus dem realen System in Paaren (\mathbf{x}_i, y_i) , wobei \mathbf{x}_i die Geschwindigkeit des Roboters und y_i die Kraft ist (angenommen, wir können die Geschwindigkeit einstellen und die Kraft mit einem Gerät messen).

Du bekommst zwei Datensätze von zwei verschiedenen Sensoren $D_a = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ und $D_b = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, in den Dateien `x_a.npy` und `y_a.npy` bzw. `x_b.npy` und `y_b.npy`.

Implementiere ein Programm in Python, das die optimalen Parameter $\boldsymbol{\theta}^*$ von der *Linear Least Squares Regression* berechnet. Du kannst als Ausgangspunkt die in der .zip-Datei enthaltene Python-Datei verwenden. Nimm an, dass dein Modell linear in den Parametern ist, $\hat{y} = \boldsymbol{\theta}^T \boldsymbol{\psi}(\mathbf{x})$, und die *feature function* ist $\boldsymbol{\psi}(\mathbf{x}) : \mathbf{x} \mapsto [1, \mathbf{x}]^T$.

Gib für jeden Datensatz die optimalen Parameter $\boldsymbol{\theta}^*$ und den *Mean Squared Error* (MSE) $L(\boldsymbol{\theta})$ an. Zeichne die gesammelten Daten und die Vorhersage der gesammelten Daten (z.B. mit der `scatter`-Funktion von Matplotlib), sowie eine Kurve mit der Vorhersage für jeden Punkt im Intervall der gesammelten Daten (z.B. mit der `plot`-Funktion von Matplotlib), d.h. wenn die gesammelten Daten \mathbf{x} zwischen -2 und 2 liegen, zeichne die Funktion zwischen -2 und 2.

Bitte lade auch deine bearbeitete Python-Datei hoch. Ändere den Namen der Datei nicht. Füge einfach deine Informationen (Name, Immatrikulationsnummer, Gruppennummer) am Anfang der Datei hinzu.

Lösungsvorschlag: (3.0 Punkt(e)) Siehe die Lösung in der Python Datei.

- e) Die Daten von d) wurden mit zwei verschiedenen Sensoren erstellt, einem sehr billigen, daher haben die Messungen ein hohes Rauschen, und einem sehr teuren, und daher haben die Messungen weniger Rauschen. Welcher Datensatz wurde mit welchem Sensor erzeugt?

Lösungsvorschlag: (1.0 Punkt(e)) Wir können die *Mean Squared Error* (MSE) vergleichen. Die Daten mit der niedrigsten MSE stammen von dem Sensor, **der Ergebnisse mit weniger Rauschen liefert.**

- f) Die in d) angenommene lineare Beziehung ist vielleicht nicht der richtige Weg, um das Kraft-Geschwindigkeits-Verhältnis zu modellieren. Hier haben wir mehr Daten in einer anderen Art von Gelände gesammelt. Wir wissen nicht mehr, ob es eine lineare Beziehung zwischen Kraft und Geschwindigkeit gibt, so dass wir keine lineare *identity feature function* mehr verwenden können. Deshalb werden wir nun verschiedene Typen von *feature functions* verwenden, um zu sehen, wie sich die Anpassung der Daten verhält. Wir werden die folgenden *feature functions* verwenden

Linear features: $\psi(\mathbf{x}) = [1, \mathbf{x}]^T$

Quadratic features: $\psi(\mathbf{x}) = [1, \mathbf{x}, \mathbf{x}^2]^T$

Cubic features: $\psi(\mathbf{x}) = [1, \mathbf{x}, \mathbf{x}^2, \mathbf{x}^3]^T$

Quartic features: $\psi(\mathbf{x}) = [1, \mathbf{x}, \mathbf{x}^2, \mathbf{x}^3, \mathbf{x}^4]^T$

Tenth features: $\psi(\mathbf{x}) = [1, \mathbf{x}, \mathbf{x}^2, \mathbf{x}^3, \dots, \mathbf{x}^{10}]^T$

Twentieth features: $\psi(\mathbf{x}) = [1, \mathbf{x}, \mathbf{x}^2, \mathbf{x}^3, \dots, \mathbf{x}^{20}]^T$

Wichtig: die *features* korrelieren nicht zwischen den Dimensionen von \mathbf{x} , d.h.

$\psi(\mathbf{x}) = [1, \mathbf{x}, \mathbf{x}^2]^T = [1, x_1, x_2, \dots, x_{d_x}, x_1^2, \dots, x_{d_x}^2]^T$, und das analog für die andere *features*.

Die Messungen des Sensors sind im Datensatz $D_c = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ gesammelt, zu finden in den Dateien `x_c.npy` und `y_c.npy`.

Implementiere die lineare Regression mit den verschiedenen *feature functions* und für den Datensatz D_c unter Verwendung der Python-Datei von d).

Gib für jeden Typ *feature function* die optimalen Parameter θ^* und den *Mean Squared Error* (MSE) $L(\theta)$ an und zeichne die gesammelten Daten sowie die Vorhersage der gesammelten Daten (z.B. mit die `scatter`-Funktion von Matplotlib), und eine Kurve mit der Vorhersage für jeden Punkt im Intervall der gesammelten Daten (z.B. mit der `plot`-Funktion von Matplotlib), d.h. wenn die gesammelten Daten \mathbf{x} zwischen -2 und 2 liegen, zeichne zwischen -2 und 2.

Was kannst du aus den Ergebnissen über das Verhältnis zwischen der Kraft und der Geschwindigkeit sagen?

Bitte lade auch deine bearbeitete Python-Datei hoch. Ändere den Namen der Datei nicht. Füge einfach deine Informationen (Name, Immatrikulationsnummer, Gruppennummer) am Anfang der Datei hinzu.

Lösungsvorschlag: (3.0 Punkt(e)) Schau die Lösung in der Python Datei.

Von der Plots schauen wir dass das Verhältnis zwischen der Kraft und der Geschwindigkeit nicht linear ist. Wenn sie linear wäre, dann würde eine lineare *feature function* gut zu den Daten passen.

- g) Wir wollen nun ein lineares Modell in einem echtem Roboter verwenden, um in Echtzeit zu arbeiten. Wir wollen das Verhältnis zwischen der Beschleunigung des Roboters und allen anderen Sensormessungen modellieren. Wir werden annehmen, dass diese Beziehung in den Parametern unseres Modells linear ist. Zu diesem Zweck sammeln wir Daten und wollen unser Modell so schnell wie möglich an den Roboter anpassen, damit es auf dem Roboter eingesetzt werden kann. Statt dass x eine einzige reelle Zahl ist, handelt es sich jetzt um einen Vektor von Sensormessungen, und wir haben insgesamt etwa 10000 Sensoren, und vielleicht möchten wir in Zukunft sogar noch mehr hinzufügen.

Nehmen wir zunächst einmal an, dass wir nur eine *feature function* wie in d) verwenden. In der Theorie können wir die optimalen Parameter θ^* in analytischer Form berechnen. Dies könnte jedoch problematisch sein, wenn wir die Ergebnisse so schnell wie möglich erhalten wollen. Warum ist das so? Was ist das Problem bei der Verwendung der analytischen Form? Um es dir leicht zu machen, kannst du bedenken, dass du nur wenige Datenpunkte hast, zum Beispiel $N = 100$.

Hinweise: Denke an die Komplexität der Berechnung von θ^ im Sinne der Big-O-Notation.*

Lösungsvorschlag: (1.0 Punkt(e)) Das Problem bei der analytischen Lösung liegt in der Berechnung der Inversen einer Matrix. Die Matrix $\psi\psi^T$ ist $D \times D$. Die Berechnung der Inversen einer Matrix kostet $\mathcal{O}(D^3)$, so dass für große D die Berechnungszeit problematisch wird. Extra: Die Invertierung einer Matrix kann auch zu numerischen Fehlern führen.