

Einführung in die mathematische Datenanalyse

Jan Heiland

FAU Erlangen-Nürnberg – Sommersemester 2022

Contents

Vorwort	5
1 Was ist Data Science?	7
1.1 Wie passiert die Datenanalyse?	7
1.2 Was sind Daten?	8
1.3 Beispiele {#01-tabellendaten}	8
1.4 Python	11
1.5 Aufgaben	11
2 Lineare Regression	15
2.1 Rauschen und Fitting	16
2.2 Ansätze für lineare Regression	17
2.3 Fehlerfunktional und Minimierung	18
2.4 Berechnung der Bestlösung	19
2.5 Beispiel	20
3 Matrix-Zerlegungen	21
3.1 QR Zerlegung	21
3.2 Singulärwertzerlegung	22
3.3 Aufgaben	24
4 Hauptkomponentenanalyse	31
4.1 Variationskoeffizienten	31
4.2 Koordinatenwechsel	33
4.3 Maximierung der Varianz in (Haupt)-Achsenrichtung	34
5 Hauptkomponentenanalyse Ctd.	37
5.1 Der IRIS Datensatz	37
5.2 Darstellung	38
5.3 Korrelationen und die Kovarianzmatrix	38
5.4 Hauptachsentransformation	39
5.5 Rekonstruktion	40
5.6 Reduktion der Daten	40
5.7 Am Beispiel der Iris Daten	41

5.8 Aufgaben	42
------------------------	----

Vorwort

Das ist ein Aufschrieb.

Korrekturen und Wünsche immer gerne als *issues* oder *pull requests* ans [github-repo](#).

Chapter 1

Was ist Data Science?

Data Science umfasst unter anderem folgende Aufgaben:

1. Strukturieren/Aufbereiten (Umgehen mit falschen, korruptierten, fehlenden, unformatierten Daten)
2. Data Exploration (Daten “verstehen”)
3. Data Analysis (quantitative Analysen, Hypothesen aufstellen)
4. Data Visualization (Hypothesen graphisch kommunizieren)
5. Modelle erzeugen/validieren (Regeln/Muster erkennen, Vorhersagen treffen) – *das* ist Machine Learning aber es gibt auch viele andere Ansätze.
6. Daten Reduktion

1.1 Wie passiert die Datenanalyse?

Mit mathematischen Methoden aus den Bereichen der

- linearen Algebra (z.B. Matrizen, Basen, lineare Gleichungssysteme)
- Statistik (z.B. Mittelwerte, Korrelationen, Verteilungen)
- Analysis (Grenzwerte, Abschätzungen)
- ...

Dabei hilft Software, z.B.,

- Excel
- **Python**
- Matlab
- R

bei der Berechnung, Automatisierung, Visualisierung.

Python ist ein de-facto Standard in Data Science und Machine Learning.

1.2 Was sind Daten?

Wie sehen Daten aus?

- Numerisch reell, z.B. Temperatur
- Numerisch diskret, z.B. Anzahl
- Ordinal: Element einer festen Menge mit expliziter Ordnung, z.B. {neuwertig, mit Gebrauchsspuren, defekt}
- Binär: Eine von zwei Möglichkeiten, z.B. Wahr/Falsch oder aktiv/inaktiv
- Kategoriell: Element einer festen Menge ohne klare Ordnung, z.B. {Säugetier, Vogel, Fisch}
- sonstige strukturierte Daten, z.B. geographische Daten, Graphen
- reiner Text, z.B. Freitext in Restaurantbewertung

Außerdem können wir noch allgemeine Eigenschaften (Qualitätsmerkmale) von Daten unterscheiden

- strukturiert
- lückenhaft
- fehlerbehaftet (*verrauscht*)
- interpretierbar
- geordnet (oder nicht zu ordnen)

1.3 Beispiele {#01-tabellendaten}

1.3.1 Tabellendaten – Mietpreise

Hier wären die Aufgaben von Data Science:

- Daten “verstehen”, Zusammenhänge zwischen Variablen aufdecken,
- visualisieren.
- Gegebenenfalls fehlende Einträge bei (z.B.) `kaltmiete` vorhersagen

Datenexploration und -analyse für einzelne Variablen 1/3 Wir betrachten eine *numerische* Variable in einem rechteckigen Datensatz, also eine *Spalte* (z.B. `kaltmiete`). Wir bezeichnen den i -ten Eintrag in dieser Spalte mit x_i , wobei $i = 1, \dots, N$ (N Anzahl der Zeilen).

Folgende *Schätzer/Metriken* können dabei helfen, diese Spalte besser zu verstehen:

- Mittelwert $\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$

	bundesland	stadt	baujahr	etage	hat_kueche	kaltmiete	wohnflaeche	zimmer
0	Nordrhein_Westfalen	Duisburg	1973	1	0	640	105	3
1	Baden_Württemberg	Ulm	1936	1	0	302	43	2
2	Nordrhein_Westfalen	Wuppertal	1986	3	0	150	67	2
3	Rheinland_Pfalz	Rhein_Lahn_Kreis	1970	3	0	315	47	3
4	Sachsen	Chemnitz	1900	2	0	315	63	2
5	Rheinland_Pfalz	Mainz	1989	1	0	1200	96	4
6	Bayern	Aschaffenburg_Kreis	1965	2	0	722	85	3
7	Berlin	Berlin	1952	0	1	706	63	2
8	Sachsen	Mittelsachsen_Kreis	1996	3	1	220	29	1
9	Brandenburg	Potsdam	1985	2	1	400	34	1

Figure 1.1: Abbildung: Tabelle von Wohnungsangeboten

	bundesland	stadt	baujahr	etage	hat_kueche	kaltmiete	wohnflaeche	zimmer
0	Nordrhein_Westfalen	Duisburg	1973	1	0	640	105	3
1	Baden_Württemberg	Ulm	1936	1	0	302	43	2
2	Nordrhein_Westfalen	Wuppertal	1986	3	0	150	67	2
3	Rheinland_Pfalz	Rhein_Lahn_Kreis	1970	3	0	315	47	3
4	Sachsen	Chemnitz	1900	2	0	315	63	2
5	Rheinland_Pfalz	Mainz	1989	1	0	1200	96	4
6	Bayern	Aschaffenburg_Kreis	1965	2	0	722	85	3
7	Berlin	Berlin	1952	0	1	706	63	2
8	Sachsen	Mittelsachsen_Kreis	1996	3	1	220	29	1
9	Brandenburg	Potsdam	1985	2	1	400	34	1

Figure 1.2: Abbildung: Eine Spalte der Tabelle

- gewichteter Mittelwert $\bar{x}_w = \frac{\sum_{i=1}^N w_i x_i}{\sum_{j=1}^N w_j}$, wobei w_i das Gewicht des i -ten Eintrages ist (z.B. eine andere Variable).
- Varianz: $s_x^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2$
- Standardabweichung $s = \sqrt{s_x^2}$.
- Median = $\frac{315+400}{2} = 357.5$.

Datenexploration und -analyse für mehrere Variablen Wir betrachten zwei Spalten $x = (x_1, \dots, x_N)$ und $y = (y_1, \dots, y_N)$.

Die Verteilung von zwei Variablen lässt sich im sogenannte **Scatter Plot** visualisieren.



Datenexploration und -analyse für mehrere Variablen Wir betrachten zwei Spalten $x = (x_1, \dots, x_N)$ und $y = (y_1, \dots, y_N)$.

- Kovarianz $s_{xy} = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$
- Korrelation $\rho_{xy} = \frac{s_{xy}}{s_x \cdot s_y} \in [-1, 1]$.
- $\rho \approx 1$: Starke positive Korrelation, wenn x groß ist, ist y auch groß.
- $\rho \approx -1$: Starke negative Korrelation, wenn x groß ist, ist y klein
- $\rho \approx 0$: Wenig/keine Korrelation.

1.3.2 COVID-19 Daten

Vergleiche die Einführung in *Mathematik für Data Science 1* vom letzten Semester.

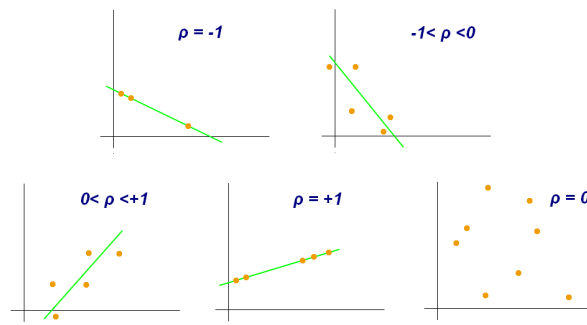


Figure 1.3: Von Kiatdd - Eigenes Werk, CC BY-SA 3.0, <https://commons.wikimedia.org/w/index.php?curid=37108966>

1.3.3 Netflix Prize

Hierbei geht es darum, ob aus bekannten Bewertungen von vielen verschiedenen Benutzern für viele verschiedene Filme abgeleitet werden kann, ob ein bestimmter Nutzer einen bestimmten Film mag (also positiv bewerten würde).

Vergleiche auch [Wikipedia:Netflix_Prize](#)

Das (Trainings-)Daten bestehen über 480189 Benutzer, die für 17770 Filme insgesamt 100480507 Bewertungen als ganze Zahlen zwischen 1 und 5 verteilen.

Ziel der Datenanalyse war es, für 2817131 “Paare” von Benutzern und Filmen, die Bewertung vorauszusagen. Neben der schieren Masse an Daten kamen noch Einschränkungen hinzu, die ein Mindestmaß an Qualität der Vorhersage sicherstellen sollten.

Das Problem ließe sich wie folgt darstellen.

Benutzer \ Film	F1	F2	...	F _n	...
B1	–	3	...	5	...
B2	3	4	...	2	...
B3	1	2	...	?	...
...	3	4	...	–	...

Gegeben viele (aber bei weitem nicht alle) Einträge in einer riesigen Tabelle. Können wir aus den Zusammenhängen bestimmte fehlende Einträge (z.B. wie findet Nutzer B3 den Film F_n) herleiten?

Die besten Lösungen für dieses Problem basieren durchweg auf *Machine Learning* Ansätzen.

1.4 Python

Die Programmiersprache `python` wird uns durchs Semester begleiten. Einfach weil sie so wichtig ist für *Data Science* aber auch weil sie (meiner Meinung nach) einfach zu erlernen und zu benutzen ist.

1.5 Aufgaben

1.5.1 Python

Bringen sie ihr `python` zum Laufen, installieren sie `numpy`, `scipy` und

```

N = 20
xmax = 2
xmin = 0

xdata = np.linspace(xmin, xmax, N)
ydata = np.exp(xdata)

plt.figure(1)
plt.plot(xdata, ydata, '.')
```

```

plt.figure(2)
plt.semilogy(xdata, ydata, '.')
```

```

plt.show()
```

1.5.2 Einheitsmatrix

Schreiben sie ein script, dass die 5x5 Einheitsmatrix auf 3 verschiedene Arten erzeugt. (Eine Art könnte die eingebaute `numpy` Funktion `eye` sein).

```

import numpy as np

idfive = np.eye(5)
print(idfive)
```

Hinweis: schauen sie sich mal an wie `numpy`'s `arrays` funktionieren.

1.5.3 Matrizen Multiplikation und Potenz

Schreiben sie ein script, das die Übungsaufgabe aus der Vorlesung (potenzieren der Matrizen M_i , $i = 1, 2, 3, 4$) löst. Zum Beispiel mit

```

import numpy as np
mone = np.array([[0.9, 0.9], [0.9, 0.9]])

mone_ptwo = mone @ mone
print(mone_ptwo)

mone_pfour = mone_ptwo @ mone_ptwo
print(mone_pfour)
```

Oder so:

```

import numpy as np
mone = np.array([[0.9, 0.9], [0.9, 0.9]])
mone_p = np.eye(2)

for k in range(16):
```

```
mone_p = mone_p @ mone
if k == 1 or k == 3 or k == 15:
    print('k=', k+1)
    print(mone_p)
```

Achtung:

- bei Matrizen kann auch `*` benutzt werden – das ist aber nicht die richtige Matrizenmultiplikation (sondern die Multiplikation eintragsweise)
- Mögliche Realisierung der Matrizenmultiplikation
 - `np.dot(A, B)` – die klassische Methode
 - `A.dot(B)` – das selbe (manchmal besser, wenn `A` etwas allgemeiner ist (zum Beispiel eine `scipy.sparse` matrix)
 - `A @ B` – convenience Notation

Chapter 2

Lineare Regression

Auch bekannt als

- *lineare Ausgleichsrechnung* oder
- *Methode der kleinsten Quadrate*.

Ein wesentlicher Aspekt von *Data Science* ist die Analyse oder das Verstehen von Daten. Allgemein gesagt, es wird versucht, aus den Daten heraus Aussagen über Trends oder Eigenschaften des Phänomens zu treffen, mit welchem die Daten im Zusammenhang stehen.

Wir kommen nochmal auf das Beispiel aus der Einführungswoche zurück, werfen eine bereits geschärften Blick darauf und gehen das mit verbesserten mathematischen Methoden an.

Gegeben seien die Fallzahlen aus der CoVID Pandemie 2020 für Bayern für den Oktober 2020.

Table 2.1: Anzahl der SARS-CoV-2 Neuinfektionen in Bayern im Oktober 2020.

Tag	1	2	3	4	5	6	7	8	9	10	11
Fälle	352	347	308	151	360	498	664	686	740	418	320

Tag	12	13	14	15	16	17	18	19	20	21
Fälle	681	691	1154	1284	127	984	573	1078	1462	2239

Tag	22	23	24	25	26	27	28	29	30	31
Fälle	2236	2119	1663	1413	2283	2717	3113	2972	3136	2615



Figure 2.1: Fallzahlen von Sars-CoV-2 in Bayern im Oktober 2020

Wieder stellen wir uns die Frage ob wir **in den Daten einen funktionalen Zusammenhang** feststellen können. Also ob wir die Datenpaare

(Tag x , Infektionen am Tag x)

die wir als

(x_i, y_i)

über eine Funktion f und die Paare

$(x, f(x))$

beschreiben (im Sinne von gut darstellen oder approximieren) können.

2.1 Rauschen und Fitting

Beim obigen Beispiel (und ganz generell bei Daten) ist davon auszugehen, dass die Daten **verrauscht** sind, also einem Trend folgen oder in einem funktionalen Zusammenhang stehen aber zufällige Abweichungen oder Fehler enthalten.

Unter diesem Gesichtspunkt ist eine Funktion, die

$$f(x_i) = y_i$$

erzwingt nicht zielführend. (Wir wollen Trends und größere Zusammenhänge erkennen und nicht kleine Fehler nachzeichnen.) Das zu strenge Anpassen an möglicherweise verrauschte Daten wird **overfitting** genannt.

Vielmehr werden wir nach einer Funktion f suchen, die die Daten näherungsweise nachstellt:

$$f(x_i) \approx y_i$$

Hierbei passen jetzt allerdings auch Funktionen, die vielleicht einfach zu handhaben sind aber die Daten kaum noch repräsentieren. Jan spricht von **underfitting**.

Eine gute Approximation besteht im Kompromiss von *nah an den Daten* aber mit wenig *overfitting*.

2.2 Ansätze für lineare Regression

Um eine solche Funktion f zu finden, trifft Jan als erstes ein paar Modellannahmen. Modellannahmen legen fest, wie das f im Allgemeinen aussehen soll und versuchen dabei

1. die Bestimmung von f zu ermöglichen
2. zu garantieren, dass f auch die gewollten Aussagen liefert
3. und sicherzustellen, dass f zum Problem passt.

Jan bemerke, dass die ersten beiden Annahmen im Spannungsverhältnis zur dritten stehen.

Lineare Regression besteht darin, dass die Funktion als Linearkombination

$$f_w(x) = \sum_{j=1}^n w_j b_j(x)$$

von Basisfunktionen geschrieben wird und dann die *Koeffizienten* w_i so bestimmt werden, dass f die Daten bestmöglich annähert.

Jan bemerke, dass *bestmöglich* wieder *overfitting* bedeuten kann aber auch, bei schlechter Wahl der Basis, wenig aussagekräftig sein kann. Der gute Kompromiss liegt also jetzt in der Wahl der passenden Basisfunktionen und deren Anzahl. (Mehr Basisfunktionen bedeutet möglicherweise bessere Approximation aber auch die Gefahr von *overfitting*.)

Typische Wahlen für die Basis $\{b_1, b_2, \dots, b_n\}$ sind

- Polynome: $\{1, x, x^2, \dots, x^{n-1}\}$ – für $n = 2$ ist der Ansatz *eine Gerade*

- Trigonometrische Funktionen: $\{1, \cos(x), \sin(x), \cos(2x), \sin(2x), \dots\}$
- Splines – Polynome, die abschnittsweise definiert werden
- Wavelets – Verallgemeinerungen von trigonometrischen Funktionen

2.3 Fehlerfunktional und Minimierung

Wir setzen nun also an

$$f_w(x) = \sum_{j=1}^n w_j b_j(x)$$

und wollen damit $y_i \approx f_w(x_i)$ *bestmöglich* erreichen (indem wir die Koeffizienten (w_1, \dots, w_n) *optimal* wählen. Bestmöglich und optimal spezifizieren wir über den Mittelwert der quadratischen Abweichungen in der Approximation über alle Datenpunkte

$$\frac{1}{N} \sum_{i=1}^N (y_i - f_w(x_i))^2$$

Ein paar Bemerkungen

- jetzt müssen wir die w_i 's bestimmen so dass dieser Fehlerterm minimal wird
- das optimale w ist unabhängig von einer Skalierung des Fehlerterms
- deswegen schreiben wir gerne einfach $\frac{1}{2} \sum_{i=1}^N (y_i - f_w(x_i))^2$ als das Zielfunktional, das es zu minimieren gilt.

Wie finden wir jetzt die w_i 's? Zunächst gilt, dass

$$f_w(x_i) = \sum_{j=1}^n w_j b_j(x_i) = \begin{bmatrix} b_1(x_i) & b_2(x_i) & \dots & b_n(x_i) \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

und wenn wir alle $f_w(x_i)$, $i = 1, \dots, N$ übereinander in einen Vektor schreiben, dass

$$f_w(\mathbf{x}) := \begin{bmatrix} f_w(x_1) \\ \vdots \\ f_w(x_N) \end{bmatrix} = \begin{bmatrix} b_1(x_1) & \dots & b_n(x_1) \\ \vdots & \ddots & \vdots \\ b_1(x_N) & \dots & b_n(x_N) \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} =: \Phi(\mathbf{x})w$$

Damit, mit \mathbf{y} als den Vektor aller y_i 's, und mit der Definition der Vektornorm, können wir unser Minimierungsproblem schreiben als

$$\frac{1}{2} \sum_{i=1}^N (y_i - f_w(x_i))^2 = \frac{1}{2} \|\mathbf{y} - \Phi(\mathbf{x})w\|^2 \rightarrow \min.$$

Wir bemerken, dass

- das Fehlerfunktional immer größer und bestenfalls gleich 0 ist
- falls das lineare Gleichungssystem $\Phi(\mathbf{x})w = \mathbf{y}$ eine Lösung w hat, ist das auch eine Lösung unserer Minimierung
- im typischen Falle aber ist allerdings $N \gg n$ und das System überbestimmt ($n = N$ würde ein overfitting bedeuten...) sodass wir keine Lösung des linearen Gleichungssystems erwarten können.
- Das Minimierungsproblems selbst hat allerdings immer eine Lösung.

2.4 Berechnung der Bestlösung

Wir suchen also ein Minimum der Funktion (mit Φ , \mathbf{x} , \mathbf{y} gegeben)

$$\begin{aligned} w \mapsto \frac{1}{2} \|\mathbf{y} - \Phi(\mathbf{x})w\|^2 &= \frac{1}{2} (\mathbf{y} - \Phi(\mathbf{x})w)^T (\mathbf{y} - \Phi(\mathbf{x})w) \\ &= \frac{1}{2} [\mathbf{y}^T \mathbf{y} - \mathbf{y}^T \Phi(\mathbf{x})w - w^T \Phi(\mathbf{x})^T \mathbf{y} + w^T \Phi(\mathbf{x})^T \Phi(\mathbf{x})w] \\ &= \frac{1}{2} [\mathbf{y}^T \mathbf{y} - 2w^T \Phi(\mathbf{x})^T \mathbf{y} + w^T \Phi(\mathbf{x})^T \Phi(\mathbf{x})w] \end{aligned}$$

wobei wir die Definition der Norm $\|v\|^2 = v^T v$ und die Eigenschaft, dass für die skalare Größe $w^T \Phi(\mathbf{x})^T \mathbf{y} = [w^T \Phi(\mathbf{x})^T \mathbf{y}]^T = \mathbf{y}^T \Phi(\mathbf{x})w$ gilt, ausgenutzt haben.

Wären w und \mathbf{y} keine Vektoren sondern einfach reelle Zahlen, wäre das hier eine Parabelgleichung $aw^2 + bw + c$ mit $a > 0$, die immer eine Minimalstelle hat.

Tatsächlich gilt hier alles ganz analog. Insbesondere ist $\Phi(\mathbf{x})^T \Phi(\mathbf{x})$ in der Regel "größer 0" (was heißt das wohl bei quadratischen Matrizen?). Und mittels "Nullsetzen" der ersten Ableitung können wir das Minimum bestimmen. In diesem Fall ist die erste Ableitung (nach w)

$$\nabla_w \left(\frac{1}{2} \|\mathbf{y} - \Phi(\mathbf{x})w\|^2 \right) = \Phi(\mathbf{x})^T \Phi(\mathbf{x})w - \Phi(\mathbf{x})^T \mathbf{y},$$

(den *Gradienten* ∇_w als Ableitung von Funktionen mit mehreren Veränderlichen werden wir noch genauer behandeln) was uns als Lösung, die Lösung des linearen Gleichungssystems

$$\Phi(\mathbf{x})^T \Phi(\mathbf{x})w = \Phi(\mathbf{x})^T \mathbf{y}$$

definiert.

Letzte Frage: Wann hat dieses Gleichungssystems eine eindeutige Lösung? Mit $N > n$ (also $\Phi(\mathbf{x})$ hat mehr Zeilen als Spalten) gelten die Äquivalenzen:

- $\Phi(\mathbf{x})^T \Phi(\mathbf{x})w = \Phi(\mathbf{x})^T \mathbf{y}$ hat eine eindeutige Lösung
- die Matrix $\Phi(\mathbf{x})^T \Phi(\mathbf{x})$ ist regulär
- die Spalten von $\Phi(\mathbf{x})$ sind linear unabhängig
- die Vektoren $b_i(\mathbf{x})$ sind linear unabhängig.

Praktischerweise tritt genau diese Situation im Allgemeinen ein.

- $N > n$ (mehr Datenpunkte als Parameter)
- b_i 's werden als *linear unabhängig* (im Sinne ihres Funktionenraums) gewählt, was die lineare unabhängigkeit der $b_i(\mathbf{x})$ impliziert.

2.5 Beispiel

Unsere Covid-Zahlen “mit einer Geraden angenähert”:

- $f_w(x) = w_1 + w_2x$ – das heißt $n = 2$ und Basisfunktionen $b_1(x) \equiv 1$ und $b_2(x) = x$
- $\mathbf{x} = (1, 2, 3, \dots, 31)$ – die Tage im Februar, das heißt $N = 31$
- $\mathbf{y} = (352, 347, \dots, 2615)$ – die Fallzahlen

Wir bekommen

$$\Phi(\mathbf{x}) = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ \vdots & \vdots \\ 1 & 31 \end{bmatrix}$$

(die Spalten sind linear unabhängig) und müssen “nur” das 2x2 System

$$\begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & 2 & 3 & \dots & 31 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ \vdots & \vdots \\ 1 & 31 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & 2 & 3 & \dots & 31 \end{bmatrix} \begin{bmatrix} 352 \\ 347 \\ 308 \\ \vdots \\ 2615 \end{bmatrix}$$

lösen um die Approximation f_w zu bestimmen.

Und noch als letzte Bemerkung. Egal wie die Basisfunktionen b_i gewählt werden, die Parameterabhängigkeit von w ist immer linear. Deswegen der Name **lineare Ausgleichsrechnung**.

Chapter 3

Matrix-Zerlegungen

Die Lösung w des Problems der *linearen Ausgleichsrechnung* war entweder als Lösung eines Optimierungsproblems

$$\min_w \|Aw - y\|^2$$

oder als Lösung des linearen Gleichungssystems

$$A^T Aw = y$$

gegeben. Hierbei steht nun $A \in \mathbb{R}^{N \times n}$ für die Matrix $\Phi(\mathbf{x})$ der Daten und Basisfunktionen. Wir hatten uns überlegt, dass in den meisten Fällen

- die Matrix mehr Zeilen als Spalten hat ($N > n$) und
- die Spalten linear unabhängig sind.

3.1 QR Zerlegung

Wir betrachten nochmal das Optimierungsproblem $\min_w \|Aw - y\|^2$. Gäbe es eine Lösung des Systems $Aw = y$, wäre das sofort eine Lösung des Optimierungsproblems. Da A aber mehr Zeilen als Spalten hat, ist das System $Aw = y$ überbestimmt und eine Lösung in der Regel nicht gegeben.

Die Überlegung ist nun, die Gleichung $Aw = y$ *so gut wie möglich* zu erfüllen, indem wir die relevanten Gleichungen identifizieren und wenigstens diese lösen. Ein systematischer (und wie wir später sehen werden auch zum Optimierungsproblem passender) Zugang bietet die QR Zerlegung.

Theorem 3.1 (QR Zerlegung). *Sei $A \in \mathbb{R}^{m \times n}$, $m > n$. Dann existiert eine orthonormale Matrix $Q \in \mathbb{R}^{m \times m}$ und eine obere Dreiecksmatrix $\hat{R} \in \mathbb{R}^{n \times n}$ derart dass*

$$A = QR =: Q \begin{bmatrix} \hat{R} \\ 0 \end{bmatrix}.$$

Hat A vollen (Spalten)Rang, dann ist \hat{R} invertierbar.

Hier heißt *orthonormale Matrix* Q , dass die Spalten von Q paarweise orthogonal sind. Insbesondere gilt

$$Q^T Q = I.$$

Für unser zu lösendes Problem ergibt sich dadurch die Umformung

$$Aw = y \Leftrightarrow QRw = y \Leftrightarrow Q^T QRw = Q^T y \Leftrightarrow Rw = Q^T y$$

oder auch

$$\begin{bmatrix} \hat{R} \\ 0 \end{bmatrix} w = Q^T y = \begin{bmatrix} Q_1^T \\ Q_2^T \end{bmatrix} y$$

(wobei wir die $Q_1 \in \mathbb{R}^{m \times n}$ die Matrix der ersten n Spalten von Q ist) und als *Kompromiss* der Vorschlag, das Teilsystem

$$\hat{R}w = Q_1^T y$$

nach w zu lösen und in Kauf zu nehmen, dass der Rest, nämlich das $Q_2^T y$, nicht notwendigerweise gleich null ist.

Wir halten zunächst mal fest, dass

- Obwohl Q eine reguläre Matrix ist, bedarf der Übergang von $Aw = y$ zu $Q^T Aw = Q^T y$ einer genaueren Analyse.
- Wir bemerken, dass für eine hypothetische komplette Lösung $Aw = y$, diese Transformation keine Rolle spielt.
- Für die Kompromisslösung jedoch schon, weil beispielsweise verschiedene Konstruktionen eines invertierbaren Teils, verschiedene Residuen bedeuten und somit Optimalität im Sinne von $\min_w \|Aw - y\|^2$ nicht garantiert ist.

Allerdings, wie Sie als Übungsaufgabe nachweisen werden, löst dieser Ansatz tatsächlich das Optimierungsproblem.

3.2 Singulärwertzerlegung

Eine weitere Matrix Zerlegung, die eng mit der Lösung von Optimierungsproblemen oder überbestimmten Gleichungssystemen zusammenhängt ist die *Singulärwertzerlegung* (SVD – von englisch: *Singular Value Decomposition*).

Theorem 3.2 (Singulärwertzerlegung). *Sei $A \in \mathbb{C}^{m \times n}$, $m \geq n$. Dann existieren orthogonale Matrizen $U \in \mathbb{C}^{m \times m}$ und $V \in \mathbb{C}^{n \times n}$ und eine Matrix*

$\Sigma \in \mathbb{R}^{m \times n}$ der Form

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \ddots & \vdots \\ 0 & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \sigma_n \\ 0 & 0 & \dots & 0 \\ \vdots & \ddots & & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix}$$

mit reellen sogenannten Singulärwerten

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$$

sodass gilt

$$A = U \Sigma V^*$$

wobei gilt $V^* = \overline{V^T}$ (transponiert und komplex konjugiert).

Ein paar Bemerkungen.

- Ist A reell, können auch U und V reell gewählt werden.
- Die Annahme $m \geq n$ war nur nötig um für die Matrix Σ keine Fallunterscheidung zu machen. (Für $m \leq n$ “steht der Nullblock rechts von den Singulärwerten”). Insbesondere gilt $A^* = V \Sigma U^*$ ist eine SVD von A^* .
- Eine Illustration der Zerlegung ist in Abbildung 3.1 zu sehen.

Wir machen einige Überlegungen im Hinblick auf große Matrizen. Sei dazu $m > n$, $A \in \mathbb{C}^{m \times n}$ und $A = U \Sigma V^*$ eine SVD wie in Theorem @ref{thm:SVD}. Sei nun

$$U = [U_1 \quad U_2]$$

partitioniert sodass U_1 die ersten n Spalten von U enthält.

Dann gilt (nach der Matrix-Multiplikations Regel *Zeile mal Spalte* die Teile U_2 und V_2 immer mit dem Nullblock in Σ multipliziert werden) dass

$$A = U \Sigma V = [U_1 \quad U_2] \begin{bmatrix} \hat{\Sigma} \\ 0 \end{bmatrix} V^* = U_1 \hat{\Sigma} V^*$$

Es genügt also nur die ersten m Spalten von U zu berechnen. Das ist die sogenannte **slim SVD**.

Hat, darüberhinaus, die Matrix A keinen vollen Rang, also $\text{Rg}(A) = r < n$, dann:

- ist $\sigma_i = 0$, für alle $i = r+1, \dots, n$, (wir erinnern uns, dass die Singulärwerte nach Größe sortiert sind)
- die Matrix $\hat{\Sigma}$ hat $n - r$ Nullzeilen
- für die Zerlegung sind nur die ersten r Spalten von U und V relevant – die sogenannte **Kompakte SVD**.

In der Datenapproximation ist außerdem die **truncated SVD** von Interesse. Dazu sei $\hat{r} < r$ ein beliebig gewählter Index. Dann werden alle Singulärwerte, $\sigma_i = 0$, für alle $i = \hat{r} + 1, \dots, n$, abgeschnitten – das heißt null gesetzt und die entsprechende *kompakte SVD* berechnet.

Hier gilt nun nicht mehr die Gleichheit in der Zerlegung. Vielmehr gilt

$$A \approx A_{\hat{r}}$$

wobei $A_{\hat{r}}$ aus der *truncated SVD* von A erzeugt wurde. Allerdings ist diese Approximation von A durch optimal in dem Sinne, dass es keine Matrix vom Rang $\hat{r} \geq r = \text{Rg}(A)$ gibt, die A (in der *induzierten* euklidischen Norm¹) besser approximiert. Es gilt

$$\min_{B \in \mathbb{C}^{m \times n}, \text{Rg}(B) = \hat{r}} \|A - B\|_2 = \|A - A_{\hat{r}}\|_2 = \sigma_{\hat{r}+1};$$

vgl. Bollhoefer/Mehrmann Satz 14.15.

Zum Abschluss noch der Zusammenhang zum Optimierungsproblem. Ist $A = U\Sigma V^*$ “SV-zerlegt”, dann gilt

$$A^*Aw = V\Sigma^*U^*U\Sigma V^*w = V\hat{\Sigma}^2V^*$$

und damit

$$A^*Aw = A^*y \quad \Leftrightarrow \quad V\hat{\Sigma}^2V^* = V\Sigma^*U^*y \quad \Leftrightarrow \quad w = V(\Sigma^+)^*U^*y$$

wobei

$$\Sigma^+ = \begin{bmatrix} \hat{\Sigma}^{-1} \\ 0_{m-n \times n} \end{bmatrix}$$

aus $\Sigma\hat{\Sigma}^{-1}\hat{\Sigma}^{-1}$ herrührt.

Bemerkung: Σ^+ kann auch definiert werden, wenn $\hat{\Sigma}$ nicht invertierbar ist (weil manche Diagonaleinträge null sind). Dann wird $\hat{\Sigma}^+$ betrachtet, bei welcher nur die $\sigma_i > 0$ invertiert werden und die anderen $\sigma_i = 0$ belassen werden. Das definiert eine sogenannte *verallgemeinerte Inverse* und löst auch das Optimierungsproblem falls A keinen vollen Rang hat.

3.3 Aufgaben

Erklärung: (T) heißt theoretische Aufgabe, (P) heißt programmieren.

3.3.1 Norm und Orthogonale Transformation (T)

Sei $Q \in \mathbb{R}^{n \times n}$ eine orthogonale Matrix und sei $y \in \mathbb{R}^n$. Zeigen Sie, dass

$$\|y\|^2 = \|Qy\|^2$$

gilt. (2 Punkte)

¹Auf Matrixnormen kommen wir noch in der Vorlesung zu sprechen.



Figure 3.1: Illustration der SVD. Bitte beachten der * bedeutet hier transponiert und komplex konjugiert. By Cmglee - Own work, CC BY-SA 4.0, <https://commons.wikimedia.org/w/index.php?curid=67853297>

3.3.2 Kleinste Quadrate und Mittelwert

Zeigen sie, dass der *kleinste Quadrate* Ansatz zur Approximation einer Datenwolke

$$(x_i, y_i), \quad i = 1, 2, \dots, N,$$

mittels einer konstanten Funktion $f(x) = w_1$ auf w_1 auf den Mittelwert der y_i führt. (6 Punkte)

3.3.3 QR Zerlegung und Kleinstes Quadrate Problem (T)

Sei $A \in \mathbb{R}^{m,n}$, $m > n$, A hat vollen Rang und sei

$$\begin{bmatrix} Q_1 & Q_2 \end{bmatrix} \begin{bmatrix} \hat{R} \\ 0 \end{bmatrix} = A$$

eine QR-Zerlegung von A . Zeigen sie, dass die Lösung von

$$\hat{R}w = Q_1^T y$$

ein kritischer Punkt (d.h. der Gradient ∇_w verschwindet) von

$$w \mapsto \frac{1}{2} \|Aw - y\|^2$$

ist. **Hinweis:** Die Formel für den Gradienten wurde in der Vorlesung 02 hergeleitet. (6 Punkte)

3.3.4 Eigenwerte Symmetrischer Matrizen (T)

Zeigen Sie, dass Eigenwerte symmetrischer reeller Matrizen $A \in \mathbb{R}^{n \times n}$ immer reell sind. (3 Punkte)

3.3.5 Singulärwertzerlegung und Eigenwerte (T)

Zeigen Sie, dass die quadrierten Singulärwerte einer Matrix $A \in \mathbb{R}^{m \times n}$, $m > n$, genau die Eigenwerte der Matrix $A^T A$ sind und in welcher Beziehung sie mit den Eigenwerten von AA^T stehen. **Hinweis:** hier ist “ $m > n$ ” wichtig. (6 Punkte)

3.3.6 Python – Laden und Speichern von Arrays

Speichern sie die Covid-Daten aus obiger Tabelle zur späteren Verwendung als ein `numpy.array`. Beispielsweise so:

```
import numpy as np
import matplotlib.pyplot as plt

days = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11,
         12, 13, 14, 15, 16, 17, 18, 19, 20, 21,
         22, 23, 24, 25, 26, 27, 28, 29, 30, 31]
case = [352, 347, 308, 151, 360, 498, 664, 686, 740, 418, 320,
        681, 691, 1154, 1284, 127, 984, 573, 1078, 1462, 2239,
        2236, 2119, 1663, 1413, 2283,
        2717, 3113, 2972, 3136, 2615]

data = np.vstack([days, case])
print('Shape of data: ', data.shape)

datafilestr = 'coviddata.npy'
np.save(datafilestr, data)

lddata = np.load(datafilestr)

plt.figure(1)
plt.plot(lddata[0, :], lddata[1, :], 's', label='Cases/Day')
plt.title('Covid Faelle in Bayern im Oktober 2020')
plt.legend()
plt.show()
```

3.3.7 Lineare Regression für Covid Daten (P)

Führen sie auf den Covid-daten eine lineare Regression zum Fitten

- einer konstanten Funktion $f(x) = c$

- einer linearen Funktion $f(x) = ax + b$
- einer quadratischen Funktion $f(x) = w_1 + w_2x + w_3x^2$

durch. Berechnen Sie die mittlere quadratische Abweichung $\frac{1}{N} \sum_{i=1}^N \|y_i - f(x_i)\|^2$ für alle drei Approximationen und plotten Sie die Covid-Daten zusammen mit der aus der Regression erhaltenen Funktion. Beispielsweise so:

```
import numpy as np
import matplotlib.pyplot as plt

datafilestr = 'coviddata.npy'
cvdata = np.load(datafilestr)

# ## Definieren der Basis Funktion(en)

def b_zero(x):
    ''' eine konstante Funktion '''
    return 1.

ybf = cvdata[1, :] # die y-werte
xbf = cvdata[0, :] # die x-werte

N = xbf.size # Anzahl Datenpunkte
ybf = ybf.reshape((N, 1)) # reshape als Spaltenvektor

bzx = [b_zero(x) for x in xbf] # eine Liste mit Funktionswerten
bzx = np.array(bzx).reshape((N, 1)) # ein Spalten Vektor

Phix = bzx # hier nur eine Spalte

# Das LGS AtA w = At y
rhs = Phix.T @ ybf
AtA = Phix.T @ Phix

# ACHTUNG: das hier geht nur weil AtA keine Matrix ist in diesem Fall
w = 1./AtA * rhs
# ACHTUNG: das hier ging nur weil AtA keine Matrix ist in diesem Fall

def get_regfunc(weights, basfunlist=[]):
    ''' Eine Funktion, die eine Funktion erzeugt

    Eingang: die Gewichte, eine Liste von Basisfunktionen
    Ausgang: die entsprechende Funktion `f = w_1b_1 + ...`
```

```

'''
def regfunc(x):
    fval = 0
    for kkk, basfun in enumerate(basfunlist):
        fval = fval + weights[kkk]*basfun(x)
    return fval

return regfunc

const_regfunc = get_regfunc([w], [b_zero])

const_approx = [const_regfunc(x) for x in xbf]
const_approx = np.array(const_approx).reshape((N, 1))

plt.figure(1)
plt.plot(cvdata[0, :], cvdata[1, :], 's', label='Cases/Day')
plt.plot(cvdata[0, :], const_approx, '-', label='constant fit')
plt.legend()
plt.show()

```

Hinweise:

- `liste = [f(x) for x in iterable]` ist sehr bequem um Vektoren von Funktionswerten zu erzeugen aber nicht sehr *pythonesque* (und auch im Zweifel nicht effizient). Besser ist es Funktionen zu schreiben, die *vektorisert* sind. Zum Beispiel können die meisten *built-in* Funktionen wie `np.exp` ein `array` als Eingang direkt in ein `array` der Funktionswerte umsetzen.
- *Eine Funktion, die eine Funktion erzeugt* finde ich sehr hilfreich für viele Anwendungen (ist aber manchmal nicht so gut nachvollziehbar).

3.3.8 Truncated SVD (P+T)

1. **(P)** Berechnen und plotten sie die Singulärwerte einer 4000×1000 Matrix mit zufälligen Einträgen und die einer Matrix mit “echten” Daten (hier Simulationsdaten einer Stroemungssimulation)². Berechnen sie den Fehler der *truncated SVD* $\|A - A_{\hat{r}}\|$ für $\hat{r} = 10, 20, 40$ für beide Matrizen.
2. **(T)** Was lässt sich bezüglich einer Kompression der Daten mittels SVD für die beiden Matrizen sagen. (Vergleichen sie die plots der Singulärwerte und beziehen sie sich auf die gegebene Formel für die Differenz).
3. **(P+T)** Für die “echten” Daten: Speichern sie die Faktoren der bei $\hat{r} = 40$ abgeschnittenen SVD und vergleichen Sie den Speicherbedarf der Faktoren und der eigentlichen Matrix.

²[Download bitte hier](#) – Achtung das sind 370MB

Beispielcode:

```
import numpy as np
import scipy.linalg as spla
import matplotlib.pyplot as plt

randmat = np.random.randn(4000, 1000)

rndU, rndS, rndV = spla.svd(randmat)

print('U-dims: ', rndU.shape)
print('V-dims: ', rndV.shape)
print('S-dims: ', rndS.shape)

plt.figure(1)
plt.semilogy(rndS, '.', label='Singulaerwerte (random Matrix)')

realdatamat = np.load('velfielddata.npy')

# # Das hier ist eine aufwaendige Operation
rlU, rlS, rlV = spla.svd(realdatamat, full_matrices=False)
# # auf keinen Fall `full_matrices=False` vergessen

print('U-dims: ', rlU.shape)
print('V-dims: ', rlV.shape)
print('S-dims: ', rlS.shape)

plt.figure(1)
plt.semilogy(rlS, '.', label='Singulaerwerte (Daten Matrix)')

plt.legend()
plt.show()
```

Hinweise:

- Es gibt viele verschiedene Normen für Vektoren und Matrizen. Sie dürfen einfach mit `np.linalg.norm` arbeiten. Gerne aber mal in die Dokumentation schauen *welche* Norm berechnet wird.
- Die (T) Abschnitte hier bitte mit den anderen (T) Aufgaben oder als Bildschirmausgabe im Programm.

Chapter 4

Hauptkomponentenanalyse

Während in den vorherigen Kapiteln der Versuch war, einen funktionalen Zusammenhang in Daten zu bekommen, geht es jetzt darum, statistische Eigenschaften aus Daten zu extrahieren. Wir werden sehen, dass das auch nur ein anderer Blickwinkel auf das gleiche Problem *Wie können wir die Daten verstehen?* ist und auch die SVD wieder treffen.

Weil ich den Ansatz gerne *ad hoc* also am Problem entlang motivieren und einführen will, vorweg schon mal die bevorstehenden Schritte

1. Zentrierung/Skalierung der Daten
2. Berechnung der Varianzen im Standard Koordinatensystem
3. Überlegung, dass Daten in einem anderen Koordinatensystem eventuell besser dargestellt werden im Sinne
4. Berechnung eines optimalen Koordinatenvektors mittels SVD.

Wir nehmen noch einmal die Covid-Daten her, vergessen kurz, dass es sich um eine Zeitreihe handelt und betrachten sie als Datenpunkte (x_i, y_i) , $i = 1, \dots, N$, im zweidimensionalen Raum mit Koordinaten x und y .

Als erstes werden die Daten **zentriert** indem in jeder Komponente der Mittelwert

$$x_c = \frac{1}{N} \sum_{i=1}^N x_i, \quad y_c = \frac{1}{N} \sum_{i=1}^N y_i.$$

abgezogen wird und dann noch mit dem inversen des Mittelwerts skaliert.

Also, die Daten werden durch $(\frac{x_i - \bar{x}}{\bar{x}}, \frac{y_i - \bar{y}}{\bar{y}})$ ersetzt.

4.1 Variationskoeffizienten

Als nächstes kann Jan sich fragen, wie gut die Daten durch ihren Mittelwert beschrieben werden und die Varianzen berechnen, die für zentrierte Daten so



Figure 4.1: Fallzahlen von Sars-CoV-2 in Bayern im Oktober 2020 – zentriert

aussehen

$$s_x^2 = \frac{1}{N-1} \sum_{i=1}^N x_i^2, \quad s_y^2 = \frac{1}{N-1} \sum_{i=1}^N y_i^2.$$

Im gegebenen Fall bekommen wir

$$s_x^2 \approx 0.32 \quad s_y^2 \approx 0.57$$

und schließen daraus, dass in y Richtung *viel passiert* und in x Richtung *nicht ganz so viel*. Das ist jeder Hinsicht nicht befriedigend, wir können weder

- Redundanzen ausmachen (eine Dimension der Daten vielleicht weniger wichtig?) noch
- dominierende Richtungen feststellen (obwohl dem Bild nach so eine offenbar existiert)

und müssen konstatieren, dass die Repräsentation der Daten im (x, y) Koordinatensystem nicht optimal ist.

Die Frage ist also, ob es ein Koordinatensystem gibt, dass die Daten besser darstellt.

Ein Koordinatensystem ist nichts anderes als eine Basis. Und die Koordinaten eines Datenpunktes sind die Komponenten des entsprechenden Vektors in dieser Basis. Typischerweise sind Koordinatensysteme orthogonal (das heißt eine Orthogonalbasis) und häufig noch orientiert (die Basisvektoren haben eine bestimmte Reihenfolge und eine bestimmte Richtung).

4.2 Koordinatenwechsel

Sei nun also $\{b_1, b_2\} \subset \mathbb{R}^2$ eine orthogonale Basis.

Wie allgemein gebräuchlich, sagen wir *orthogonal*, meinen aber *orthonormal*. In jedem Falle soll gelten

$$b_1^T b_1 = 1, \quad b_2^T b_2 = 1, \quad b_1^T b_2 = b_2^T b_1 = 0.$$

Wir können also alle Datenpunkte $\mathbf{x}_i = \begin{bmatrix} x_i \\ y_i \end{bmatrix}$ in der neuen Basis darstellen mit eindeutig bestimmten Koeffizienten α_{i1} und α_{i2} mittels

$$\mathbf{x}_i = \alpha_{i1} b_1 + \alpha_{i2} b_2.$$

Für orthogonale Basen sind die Koeffizienten durch *testen* mit dem Basisvektor einfach zu berechnen:

$$\begin{aligned} b_1^T \mathbf{x}_i &= b_1^T (\alpha_{i1} b_1 + \alpha_{i2} b_2) = \alpha_{i1} b_1^T b_1 + \alpha_{i2} b_1^T b_2 = \alpha_{i1} \cdot 1 + \alpha_{i2} \cdot 0 = \alpha_{i1}, \\ b_2^T \mathbf{x}_i &= b_2^T (\alpha_{i1} b_1 + \alpha_{i2} b_2) = \alpha_{i1} b_2^T b_1 + \alpha_{i2} b_2^T b_2 = \alpha_{i1} \cdot 0 + \alpha_{i2} \cdot 1 = \alpha_{i2}. \end{aligned}$$

Es gilt also

$$\alpha_{i1} = b_1^T \mathbf{x}_i = b_1^T \begin{bmatrix} x_i \\ y_i \end{bmatrix}, \quad \alpha_{i2} = b_2^T \mathbf{x}_i = b_2^T \begin{bmatrix} x_i \\ y_i \end{bmatrix}.$$

Damit, können wir jeden Datenpunkt $\mathbf{x}_i = (x_i, y_i)$ in den neuen Koordinaten $(\alpha_{i1}, \alpha_{i2})$ ausdrücken.

Zunächst halten wir fest, dass auch in den neuen Koordinaten die Daten zentriert sind. Es gilt nämlich, dass

$$\frac{1}{N} \sum_{i=1}^N \alpha_{ji} = \frac{1}{N} \sum_{i=1}^N b_j^T \mathbf{x}_i = \frac{1}{N} b_j^T \sum_{i=1}^N \begin{bmatrix} x_i \\ y_i \end{bmatrix} = \frac{1}{N} b_j^T \begin{bmatrix} \sum_{i=1}^N x_i \\ \sum_{i=1}^N y_i \end{bmatrix} = b_j^T \begin{bmatrix} \frac{1}{N} \sum_{i=1}^N x_i \\ \frac{1}{N} \sum_{i=1}^N y_i \end{bmatrix} = b_j^T \begin{bmatrix} 0 \\ 0 \end{bmatrix} = 0,$$

für $j = 1, 2$.

Desweiteren gilt wegen der Orthogonalität von $B = [b_1 \ b_2] \in \mathbb{R}^{2 \times 2}$, dass

$$x_i^2 + y_i^2 = \|\mathbf{x}_i\|^2 = \|B^T \mathbf{x}_i\|^2 = \left\| \begin{bmatrix} b_1^T \\ b_2^T \end{bmatrix} \mathbf{x}_i \right\|^2 = \left\| \begin{bmatrix} b_1^T \mathbf{x}_i \\ b_2^T \mathbf{x}_i \end{bmatrix} \right\|^2 = \left\| \begin{bmatrix} \alpha_{i1} \\ \alpha_{i2} \end{bmatrix} \right\|^2 = \alpha_{i1}^2 + \alpha_{i2}^2$$

woraus wir folgern, dass in jedem orthogonalen Koordinatensystem, die Summe der beiden Varianzen die gleiche ist:

$$s_x^2 + s_y^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i^2 + y_i^2) = \frac{1}{N-1} \sum_{i=1}^N (\alpha_{i1}^2 + \alpha_{i2}^2) =: s_1^2 + s_2^2.$$

Das bedeutet, dass durch die Wahl des Koordinatensystems die Varianz als Summe nicht verändert wird. Allerdings können wir das System so wählen, dass eine der Varianzen in Achsenrichtung maximal wird (und die übrige(n) entsprechend klein).

Analog gilt für den eigentlichen Mittelwert der (nichtzentrierten) Daten, dass die Norm gleich bleibt. In der Tat, für die *neuen* Koordinaten des Mittelwerts gilt in der Norm

$$\left\| \begin{bmatrix} \alpha_{c1} \\ \alpha_{c2} \end{bmatrix} \right\| = \left\| B^T \begin{bmatrix} x_c \\ y_c \end{bmatrix} \right\| = \left\| \begin{bmatrix} x_c \\ y_c \end{bmatrix} \right\|.$$

4.3 Maximierung der Varianz in (Haupt)-Achsenrichtung

Wir wollen nun also $b_1 \in \mathbb{R}^2$, mit $\|b_1\| = 1$ so wählen, dass

$$s_1^2 = \frac{1}{N-1} \sum_{i=1}^n \alpha_{i1}^2$$

maximal wird. Mit der Matrix \mathbf{X} aller Daten

$$\mathbf{X} = \begin{bmatrix} x_1 & y_1 \\ x_2 & y_2 \\ \vdots & \vdots \\ x_N & y_N \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} \in \mathbb{R}^{N \times 2}$$

können wir die Varianz in b_1 -Richtung kompakt schreiben als

$$s_1^2 = \frac{1}{N-1} \sum_{i=1}^n \alpha_{i1}^2 = \frac{1}{N-1} \sum_{i=1}^n (b_1^T \mathbf{x}_i)^2 = \frac{1}{N-1} \sum_{i=1}^n (\mathbf{x}_i^T b_1)^2 = \frac{1}{N-1} \|\mathbf{X} b_1\|^2$$

Wir sind also ein weiteres mal bei einem Optimierungsproblem (diesmal mit Nebenbedingung) angelangt:

$$\max_{b \in \mathbb{R}^2, \|b\|=1} \|\mathbf{X} b\|^2$$

Die Lösung dieses Problems ist mit $b = v_1$ gegeben, wobei v_1 der erste (rechte) Singulärvektor von \mathbf{X} ist:

$$\mathbf{X} = U \Sigma V^T = U \Sigma \begin{bmatrix} v_1^T \\ v_2^T \end{bmatrix}.$$

Und damit rechnen wir auch direkt nach, dass im neuen Koordinatensystem $\{b_1, b_2\} = \{v_1, v_2\}$ die Varianzen s_1^2 und s_2^2 (bis auf einen Faktor von $\frac{1}{N-1}$) genau die quadrierten Singulärwerte von \mathbf{X} sind:

$$(N-1)s_1^2 = \|\mathbf{X} v_1\|^2 = \|U \Sigma \begin{bmatrix} v_1^T \\ v_2^T \end{bmatrix} v_1\|^2 = \left\| \Sigma \begin{bmatrix} v_1^T v_1 \\ v_2^T v_1 \end{bmatrix} \right\|^2 = \left\| \Sigma \begin{bmatrix} 1 \\ 0 \end{bmatrix} \right\|^2 = \sigma_1^2,$$

$$(N-1)s_2^2 = \|\mathbf{X} v_2\|^2 = \|U \Sigma \begin{bmatrix} v_1^T \\ v_2^T \end{bmatrix} v_2\|^2 = \left\| \Sigma \begin{bmatrix} v_1^T v_2 \\ v_2^T v_2 \end{bmatrix} \right\|^2 = \left\| \Sigma \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right\|^2 = \sigma_2^2$$

4.3. MAXIMIERUNG DER VARIANZ IN (HAUPT)-ACHSENRICHTUNG35

Für unser Covid Beispiel ergibt sich

$$V^T \approx \begin{bmatrix} 0.5848 & 0.8111 \\ 0.8111 & -0.5848 \end{bmatrix}$$

also

$$b_1 = v_1 = \begin{bmatrix} 0.5848 \\ 0.8111 \end{bmatrix} \quad b_2 = v_2 = \begin{bmatrix} 0.8111 \\ -0.5848 \end{bmatrix}$$

als neue Koordinatenrichtungen mit

$$s_1^2 \approx 0.85, \quad s_2^2 \approx 0.04,$$

was bereits eine deutliche Dominanz der v_1 -Richtung – genannt *Hauptachse* – zeigt.

Im Hinblick auf die nächste Vorlesung, in der wir Anwendungen und Eigenschaften der PCA untersuchen werden, noch ein Plot der Daten mit der v_1 -Richtung als Linie eingezeichnet.



Figure 4.2: Fallzahlen von Sars-CoV-2 in Bayern im Oktober 2020 – zentriert/skaliert/Hauptachse

Chapter 5

Hauptkomponentenanalyse Ctd.

Im vorherigen Kapitel hatten wir einen zweidimensionalen Datensatz betrachtet und dafür die Richtung bestimmt, in der die Varianz maximal wird. Da wir außerdem ermittelt hatten, dass die Summe der Varianz unabhängig vom Koordinatensystem ist (es muss nur ein orthogonales sein) bedeutete das gleichzeitig, dass die verbleibende Richtung die Richtung der minimalen Varianz war.

Jetzt wollen wir einen Datensatz mit mehr Merkmalen betrachten und sehen, wie die algorithmische Herangehensweise zum Verständnis beiträgt.

5.1 Der IRIS Datensatz

Die Grundlage ist der *Iris Datensatz*, der ein gern genommene Grundlage für die Illustration der Datenanalyse ist. Die Daten wurden von einem Statistiker und Biologen erhoben und beinhalten 4 verschiedene Merkmale (engl. *features*) von einer Stichprobe von insgesamt 150 Schwertlilien die 3 verschiedenen Klassen zugeordnet werden können oder sollen (Fachbegriff hier: *targets*). Im Beispiel werden die Klassen mit 0, 1, 2 codiert und beschreiben die Unterarten *setosa*, *versicolor* und *virginica* der Schwertlilien. Die Merkmale sind gemessene Länge und Breite des Kelchblattes (hier *sepal*) und des Kronblattes (hier *petal*).

Wir stellen uns 2-3 Fragen:

1. Würden eventuell 3 (oder sogar nur 2) Dimensionen reichen um den Datensatz zu beschreiben?
2. Können wir aus den Merkmalen (*features*) die Klasse (*target*) erkennen und wie machen wir gegebenenfalls die Zuordnung?

5.2 Darstellung

In höheren Dimensionen ist schon die graphische Darstellung der Daten ein Problem. Wir können aber alle möglichen 2er Kombinationen der Daten in 2D plots visualisieren.

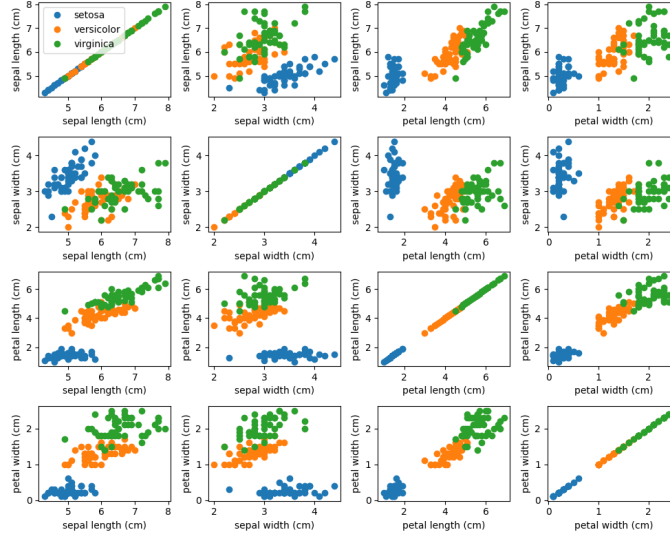


Figure 5.1: Iris Datenset 2D plots

Ein Blick auf die Diagonale zeigt schon, dass manche Merkmale besser geeignet als andere sind, um die Spezies zu unterscheiden, allerdings keine der Zweierkombinationen eine Eindeutige Diskriminierung erlaubt.

5.3 Korrelationen und die Kovarianzmatrix

Als nächstes suchen wir nach Korrelationen in den Daten in dem wir für alle Merkmalpaare die Korrelationen ausrechnen. Dafür berechnen wir die sogenannte *Kovarianzmatrix*

$$\text{Cov}(\mathbf{X}) = [\rho_{\mathbf{x}_i \mathbf{x}_j}]_{i,j=1,\dots,n} \in \mathbb{R}^{n \times n}$$

wobei n die Dimension der Daten ist und wobei

$$s_{\mathbf{x}_i \mathbf{x}_j} = \frac{1}{N-1} \sum_{k=1}^N (x_{ik} - \bar{\mathbf{x}}_i)(x_{jk} - \bar{\mathbf{x}}_j).$$

die Kovarianzen sind, die wir auch schon in der ersten Vorlesung kennengelernt haben.

Ist $\mathbf{X} \in \mathbb{R}^{N \times n}$ die Matrix mit den Datenvektoren $\mathbf{x}_i \in \mathbb{R}^n$ als Spalten **und ist der Datensatz zentriert** so erhalten wir die Kovarianzmatrix als

$$\text{Cov}(\mathbf{X}) = \frac{1}{N-1} \mathbf{X}^T \mathbf{X}$$

Wir bemerken, dass auf der Hauptdiagonalen die Varianzen in Koordinatenrichtung stehen und in den Zeilen oder Spalten ein Mass dafür, wie bspw. \mathbf{x}_i und \mathbf{x}_j korreliert sind. Große Zahlen bedeuten eine große Varianz oder eine starke Korrelation und umgekehrt. Für die Datenanalyse können wir $\text{Cov}(\mathbf{X})$ wie folgt heranziehen:

5.4 Hauptachsentransformation

Wie im vorherigen Kapitel hergeleitet, bedeutet ein Koordinatenwechsel die Multiplikation der Datenmatrix $\mathbf{X} \in \mathbb{R}^{N \times n}$ mit einer orthogonalen Matrix $V \in \mathbb{R}^{n \times n}$:

$$\tilde{\mathbf{X}} = \mathbf{X}V,$$

wobei $\tilde{\mathbf{X}}$ die Daten in den neuen Koordinaten sind (die Basisvektoren sind dann die Zeilenvektoren von V). Wir wollen nun eine Basis finden, in der

- $\text{Cov}(\mathbf{X})$ eine Diagonalmatrix ist — damit wären alle Richtungen in den Daten *unkorreliert* und könnten *unabhängig* voneinander betrachtet werden¹ –
- und in der die neuen Varianzen nach Größe geordnet sind gleichzeitig die verfügbare Varianz maximal in wenigen Richtungen konzentrieren.

Nachden den Überlegungen im vorherigen Kapitel ist die erste Hauptrichtung durch den ersten rechten Singulärvektor $v_1 \in \mathbb{R}^n$ der (*ökonomischen*) Singulärwertzerlegung

$$\mathbf{X} = U\Sigma V^* = \begin{bmatrix} u_1 & u_2 & \dots & u_n \end{bmatrix} \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_n \end{bmatrix} \begin{bmatrix} v_1^* \\ v_2^* \\ \vdots \\ v_n^* \end{bmatrix}$$

von $\mathbf{X} \in \mathbb{R}^{N \times n}$ gegeben. Die zugehörigen Koeffizienten berechnen wir mittels

$$\tilde{\mathbf{X}} = \mathbf{X}v_1.$$

¹wir dürfen aber nicht vergessen, dass Daten tüpischerweise nur eine Stichprobe von Beobachtungen eines Phänomens sind. Die Unabhängigkeit in den *features* gilt also nur für die gesammelten Daten aber in der Regel nicht für das Phänomen. Für normalverteilte Prozesse liefern die daten-basiert ermittelten Hauptrichtungen jedoch auch die Hauptrichtungen des zugrundeliegenden Phänomens

5.5 Rekonstruktion

Um zu plausibilisieren, dass die weiteren Hauptachsen durch die weiteren (rechten) Singulärvektoren gegeben sind, betrachten wir erst die *Rekonstruktion* also die Darstellung im Ausgangskoordinatensystem (mit den messbaren oder interpretierbaren Features) die durch

$$\tilde{\mathbf{X}} = \mathbf{X}v_1v_1^T$$

gegeben ist. Die letzte Formel wird vielleicht klarer, wenn Jan sich überlegt, dass für einen Datenpunkt $\mathbf{X}_j = [x_{1j} \ x_{2j} \ \dots \ x_{nj}]$ (also die j -te Zeile von \mathbf{X}), der Koeffizient für v_1 gegeben ist durch $\alpha_j = v_1^T \mathbf{X}_j^T$ und die Darstellung im Vektorraum durch

$$\tilde{\mathbf{X}}_j = (\alpha_j v_1)^T = \alpha_j v_1^T = (v_1^T \mathbf{X}_j^T) v_1^T = \mathbf{X}_j v_1 v_1^T$$

gegeben ist.

Damit können wir mit

$$\mathbf{X} - \tilde{\mathbf{X}} = \mathbf{X}(I - v_1v_1^T)$$

den Teil der Daten betrachten, der durch die Richtung v_1 nicht abgebildet wird (sozusagen den noch verbleibenden Teil). Jan kann nachrechnen, dass wiederum U und V die Singulärvektoren von $\mathbf{X} - \tilde{\mathbf{X}}$ bilden, mit jetzt v_2 als Richtung mit dem größten (verbleibenden) Singulärwert.

Das heißt, dass der k -te Singulärvektor die k -te Hauptrichtung bildet.

5.6 Reduktion der Daten

Ist die Datenmatrix mit linear abhängigen Spalten besetzt, äußert sich das in $\sigma_k = 0$ für ein $k < n$ (und alle noch folgenden Singulärwerte). Jan rechnet nach, dass dann

$$\mathbf{X} = \mathbf{X}V_{k-1}V_{k-1}^T$$

wobei V_{k-1} die Matrix der $k-1$ führenden Singulärvektoren ist und das entsprechende

$$\tilde{\mathbf{X}} = \mathbf{X}V_{k-1}$$

alle Informationen des Datensatzes in kleinerer Dimension parametrisiert.

In der Praxis sind schon allein durch Messfehler exakte lineare Abhängigkeiten sowie durch die näherungsweise Bestimmung auf dem Computer das Auftreten von $\sigma_k = 0$ quasi ausgeschlossen. Stattdessen werden Schwellwerte definiert, unterhalb derer die SVD abgeschnitten wird.

Wie oben, werden dann die Daten $\tilde{\mathbf{X}} = \mathbf{X}V_{\hat{k}}$ in den (reduzierten) Koordinaten der Hauptachsen betrachtet sowie die Rekonstruktion $\hat{\mathbf{X}} = \mathbf{X}V_{\hat{k}}V_{\hat{k}}^T$, wobei der $V_{\hat{k}}$ die Matrix der \hat{k} führenden Singulärvektoren sind (mit Singulärwerten überhalb des Schwellwertes).

5.7 Am Beispiel der Iris Daten

Wir stellen die Daten noch einmal zentriert dar:

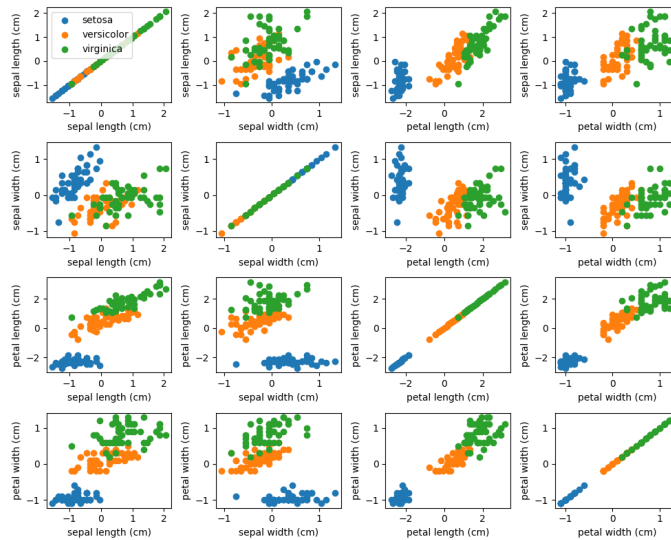


Figure 5.2: Iris Daten Zentriert

Für die Kovarianzmatrix der Iris Daten erhalten wir:

Covariance matrix:

```
[[ 0.6857 -0.0424  1.2743  0.5163]
 [-0.0424  0.19   -0.3297 -0.1216]
 [ 1.2743 -0.3297  3.1163  1.2956]
 [ 0.5163 -0.1216  1.2956  0.581  ]]
```

Für die Singulärwerte:

```
In: U, S, Vh = np.linalg.svd(data, full_matrices=False)
```

```
In: print(S)
```

```
Out: [25.1      6.0131  3.4137  1.8845]
```

Zwar ist hier kein Singulärwert nah an der Null, allerdings beträgt der Unterschied zwischen dem größten und dem kleinsten schon eine 10er Potenz, was auf eine starke Dominanz der ersten Hauptachsen hinweist.

In der Tat, plotten wir die Daten in den Koordinaten der Hauptachsen (also $\tilde{\mathbf{X}}$), zeigen Bilder der v_1 oder v_2 -Koordinaten klare Tendenzen in den Daten, während die Plots der übrigen Richtungen wie eine zufällige Punktwolke aussieht, vgl. Abbildung 5.3.

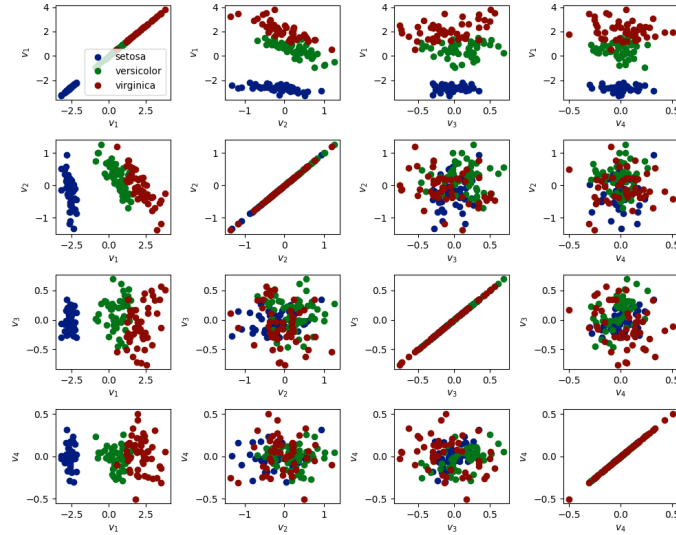


Figure 5.3: Iris Datenset 2D plots in Hauptachsenkoordinaten

Als letztes plotten wir noch $\hat{\mathbf{X}}$ für $\hat{k} = 2$. Das heißt wir reduzieren die Daten auf die v_1 und v_2 -Richtungen und betrachten die Rekonstruktion. Im Plot sehen wir, dass in gewissen Teilen die Daten gut rekonstruiert werden. Allerdings, hat $\hat{\mathbf{X}}$ nur Rang $\text{Rk } \hat{\mathbf{X}} = 2$ (warum?), sodass in den Plots notwendigerweise direkte lineare Abhängigkeiten offenbar werden.

5.8 Aufgaben

5.8.1 Kovarianzen (P+T)

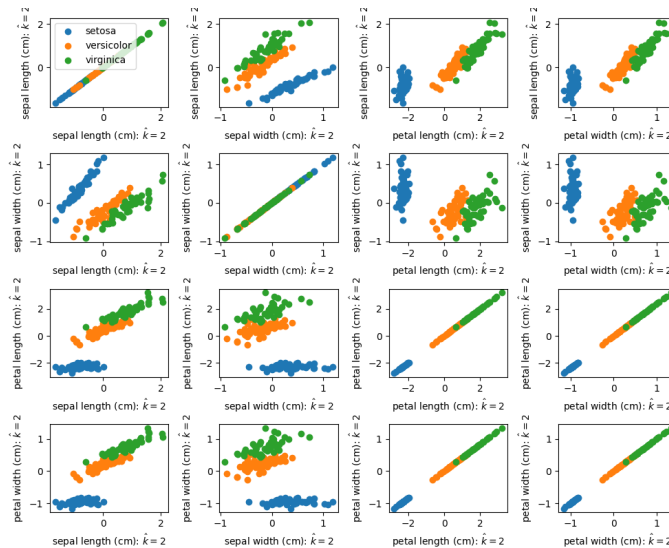
Erzeugen Sie für N aus $N1 = [10, 1000, 100000]$ zufällige Vektoren \mathbf{x} und \mathbf{y} der Länge N und berechnen Sie für den Datensatz $\mathbf{X} = [\mathbf{x}, \mathbf{x}, \mathbf{x}*\mathbf{y}, \mathbf{y}]$ jeweils die Matrix der Korrelationskoeffizienten

$$\rho_{ij} = \frac{s_{ij}}{s_i \cdot s_j}$$

wobei s_{ij} die Kovarianz der Daten \mathbf{x}_i und \mathbf{x}_j ist und s_i die Varianz von \mathbf{x}_i . Interpretieren sie die Ergebnisse.

Hinweis: Weil Zufälligkeit involviert ist, ist die Interpretation manchmal schwierig. Lassen sie das Programm öfter laufen und beobachten sie verschiedene Realisierungen der Stichproben.

```
import numpy as np
```

Figure 5.4: Iris Datenset – rekonstruiert von $\hat{k} = 2$ Hauptachsenkoordinaten

```
# N = 10

x = np.random.randn(10, 1)
y = np.random.randn(10, 1)

X = np.hstack([x, x, x*y, y])

Xcntrd = X - X.mean(axis=0)
# zentrieren der Daten

# Kovarianz Matrix ausrechnen
covX = 1/(10-1)*Xcntrd.T @ Xcntrd

# Varianzen -- stehen auf der Diagonalen im Quadrat
varvec = np.sqrt(np.diagonal(covX)).reshape(4, 1)
# Matrix mit den Kombinationen aller Varianzen
varmat = varvec @ varvec.T

matcc = covX * 1/varmat

print(f'N={10} : Matrix of Correlation Coefficients=')
print(matcc)
```

5.8.2 Iris Datensatz – Targets Plotten (P)

Laden Sie die **Iris Datensatz** (hier als json file bereitgestellt) und plotten sie petal width versus sepal length für die drei Spezies **setosa**, **versicolor**, **virginica** in drei separaten Grafiken.

```
import json

import numpy as np
import matplotlib.pyplot as plt

with open('iris-data.json', 'r') as f:
    datadict = json.load(f)

print(datadict.keys())

data = np.array(datadict['data'])
target = np.array(datadict['target'])
feature_names = datadict['feature_names']
target_names = datadict['target_names']

print('target names: ', target_names)
print('feature names: ', feature_names)

fig, axs = plt.subplots(nrows=1, ncols=3, figsize=(10, 3))

# # Ein Vektor der aus dem Target Vektor das target 0 raussucht
trgtidx_z = (target == 0)
# # Wird gleich benutzt um die Daten nach diesem target zu filtern

target_z_daten = data[trgtidx_z, :]

axs[0].plot(target_z_daten[:, 3], target_z_daten[:, 0],
            'o', label=target_names[0])

axs[0].legend()
axs[0].set_xlabel(feature_names[3])
axs[0].set_ylabel(feature_names[0])
axs[0].set_title(target_names[0])

plt.tight_layout()
plt.show()
```

Hier wurden die python Datentypen

- *list* – z.B. datenpunkte = [1, 2, 3]
- *array* – z.B. datenmatrix = np.array(datenpunkte)

- *dictionary* – z.B. `datadict = {'data': datenpunkte}`

verwendet, die alle für verschiedene Zwecke gerne benutzt werden. Z.B.

- *Liste* – als eine Sammlung von (möglicherweise total unterschiedlichen) Objekten, über die iteriert werden kann und die einfach zu erweitern ist
- *arrays* – Matrix/Vektor von Daten eines Typs, mit denen *gerechnet* werden kann
- *dictionaries* – ein *Lookup table*. Objekte können über einen Namen adressiert werden. Ich nehme sie gerne um Daten mit ihrem Namen zum Beispiel als *json* file zu speichern.

5.8.3 Iris Datensatz – 2D plots (P)

Erzeugen sie Abbildung 5.1 zum Beispiel über ein 4x4 Feld von subplots

```
fig, axs = plt.subplots(nrows=4, ncols=4, figsize=(10, 8))
```

erzeugen, in das sie mittels

```
axs[zeile, spalte].plot(xdaten, ydaten, 'o')
```

die plots fuer die einzelnen targets “eintragen” koennen. Bitte die Achsen beschriften (die *legend* ist nicht unbedingt notwendig).

5.8.4 Kovarianz (P)

Zentrieren Sie den Datensatz (zum Beispiel unter Verwendung der `numpy.mean` Funktion) und berechnen Sie die Kovarianzmatrix.

5.8.5 Hauptachsentransformation (P)

Berechnen sie Hauptachsen und stellen Sie die Daten in den Hauptachsenkoordinaten dar (wie in Abbildung 5.3).

5.8.6 Kovarianzmatrix (T)

Zeigen sie, dass für zentrierte Datensätze $\mathbf{X} \in \mathbb{R}^{N \times n}$ gilt, dass

$$\text{Cov}(X) = \frac{1}{N-1} \mathbf{X}^T \mathbf{X}.$$

5.8.7 Gesamtvarianz (T)

Zeigen sie, dass auch für $n > 2$ die Summe der Varianzen in orthogonalen Achsenrichtungen unabhängig von der Wahl des Koordinatensystems sind. (Vergleiche Kapitel 4.2 Koordinatenwechsel)