

Falls wir umgekehrt wissen, daß $(A\vec{v}) \cdot (A\vec{w}) = \vec{v} \cdot \vec{w}$ ist für alle Vektoren \vec{v}, \vec{w} , so ist auch $\vec{v} \cdot ((A^*A)\vec{w}) = \vec{v} \cdot \vec{w}$, insbesondere also

$$\vec{e}_i \cdot ((A^*A)\vec{e}_j) = \vec{e}_i \cdot \vec{e}_j = \delta_{ij}$$

für die Koordinateneinheitsvektoren.

$(A^*A)\vec{e}_j$ ist die j -te Spalte der Matrix A^*A , ihr Skalarprodukt mit \vec{e}_i also der ij -Eintrag von A^*A . Da dieser gleich δ_{ij} sein muß, ist also $A^*A = E$ und A somit orthogonal bzw. unitär. ■

Im Reellen beschreiben daher orthogonale Matrizen lineare Abbildungen, die alle Längen und Winkel respektieren. Wie wir oben gesehen haben, haben solche Matrizen entweder Determinante eins oder Determinante minus eins. Determinante minus eins tritt beispielsweise auf bei Spiegelungen, die bekanntlich im \mathbb{R}^3 nicht orientierungstreu sind. Allgemein sagt man, die lineare Abbildung zur orthogonalen Matrix A sei orientierungstreu, falls $\det A = 1$ ist. Auf die dahinter stehende Theorie der orientierten Vektorräume wollen wir nicht weiter eingehen.

h) Orthogonale Projektionen

Ist U ein r -dimensionaler Untervektorraum eines n -dimensionalen Vektorraums V , so können wir jede Basis $\{\vec{b}_1, \dots, \vec{b}_r\}$ von U ergänzen zu einer Basis $\{\vec{b}_1, \dots, \vec{b}_n\}$ von V . Der von \vec{b}_{r+1} bis \vec{b}_n erzeugte Untervektorraum $W \leq V$ hat dann die Eigenschaft, daß $U \cap W$ der Nullraum ist, während $U \cup W$ dem gesamten Vektorraum V erzeugt. Einen solchen Untervektorraum W bezeichnen wir als *Komplement* von U ; es ist natürlich, genau wie seine Basisvektoren \vec{b}_{r+1} bis \vec{b}_n , alles andere als eindeutig bestimmt.

Für EUKLIDISCHE und HERMITISCHE Vektorräume können wir allerdings jedem Untervektorraum ein wohlbestimmtes ausgezeichnetes Komplement zuordnen, das *orthogonale Komplement*.

Definition: V sei ein EUKLIDISCHER oder HERMITISCHER Vektorraum und $U \leq V$ sei ein Untervektorraum von V . Das orthogonale Komplement von U ist der Untervektorraum

$$U^\perp \stackrel{\text{def}}{=} \{ \vec{v} \in V \mid \vec{u} \cdot \vec{v} = 0 \text{ für alle } \vec{u} \in U \}.$$

Wegen der Linearität des EUKLIDISCHEN wie auch HERMITISCHEN Skalarprodukts im ersten Argument ist klar, daß U^\perp ein Untervektorraum von V ist. Außerdem ist klar, daß es reicht die Bedingung $\vec{v} \cdot \vec{u} = 0$ für die Vektoren \vec{u} aus einer Basis von U nachzurechnen, denn wenn alle diese Produkte verschwinden, verschwindet auch jedes Produkt mit einer Linearkombination solcher Vektoren. (Es stört dabei nicht, daß wir im HERMITISCHEN Fall keine Linearität im zweiten Argument haben, sondern die Koeffizienten komplex konjugieren müssen.)

Lemma: U sei ein Untervektorraum des n -dimensionalen EUKLIDISCHEN oder HERMITISCHEN Vektorraums V , und $(\vec{b}_1, \dots, \vec{b}_r)$ sei eine Orthogonalbasis von U . Ergänzt man diese zu einer Orthogonalbasis $(\vec{b}_1, \dots, \vec{b}_n)$ von V , so ist $(\vec{b}_{r+1}, \dots, \vec{b}_n)$ eine Orthogonalbasis von U^\perp . Insbesondere hat also das orthogonale Komplement eines r -dimensionalen Untervektorraum die Dimension $n - r$ und $U \cap U^\perp = \{0\}$.

Beweis: Nach dem gerade Gesagten liegt ein Vektor $\vec{v} = \lambda_1 \vec{b}_1 + \dots + \lambda_n \vec{b}_n$ aus V genau dann in U^\perp , wenn für alle $i \leq r$ gilt:

$$\vec{v} \cdot \vec{b}_i = \left(\sum_{j=1}^n \lambda_j \vec{b}_j \right) \cdot \vec{b}_i = \sum_{j=1}^n \lambda_j \vec{b}_j \cdot \vec{b}_i = \lambda_i \vec{b}_i \cdot \vec{b}_i = 0.$$

Da \vec{b}_i als Basisvektor nicht der Nullvektor sein kann, ist $\vec{b}_i \cdot \vec{b}_i \neq 0$; daher ist dies äquivalent zum Verschwinden aller λ_i mit $i \leq r$, also zur Darstellbarkeit von \vec{v} als Linearkombination der Vektoren $\vec{b}_{r+1}, \dots, \vec{b}_n$. Als Teil einer Basis sind diese linear unabhängig, also Basis ihres Erzeugnisses U^\perp . ■

Korollar: a) V sei ein endlichdimensionaler EUKLIDISCHER oder HERMITISCHER Vektorraum, und $U \leq V$ sei ein Untervektorraum. Dann läßt sich jedes Element $\vec{v} \in V$ eindeutig schreiben als $\vec{v} = \vec{u} + \vec{w}$ mit $\vec{u} \in U$ und $\vec{w} \in U^\perp$.
b) $U^{\perp\perp} = U$

Beweis: a) Wir wählen eine Orthogonalbasis $(\vec{b}_1, \dots, \vec{b}_r)$ von U und ergänzen sie zu einer Orthogonalbasis $(\vec{b}_1, \dots, \vec{b}_n)$ von V ; nach dem

Lemma ist dann $(\vec{b}_{r+1}, \dots, \vec{b}_n)$ eine Orthogonalbasis von U^\perp . Schreiben wir $\vec{v} = v_1 \vec{b}_1 + \dots + v_n \vec{b}_n$, so ist also

$$\vec{u} \stackrel{\text{def}}{=} v_1 \vec{b}_1 + \dots + v_r \vec{b}_r \in U, \quad \vec{w} \stackrel{\text{def}}{=} v_{r+1} \vec{b}_{r+1} + \dots + v_n \vec{b}_n \in U^\perp$$

und $\vec{v} = \vec{u} + \vec{w}$.

Ist $\vec{v} = \vec{x} + \vec{y}$ irgendeine Darstellung von \vec{v} als Summe zweier Vektoren $\vec{x} \in U$ und $\vec{y} \in V$, so ist

$$\vec{u} + \vec{w} = \vec{x} + \vec{y} \implies \vec{u} - \vec{x} = \vec{y} - \vec{w}.$$

In der letzteren Gleichung steht links der Vektor $\vec{u} - \vec{x} \in U$ und rechts $\vec{y} - \vec{w} \in U^\perp$; wegen $U \cap U^\perp = \{0\}$ ist also $\vec{u} = \vec{x}$ und $\vec{w} = \vec{y}$, was die Eindeutigkeit dieser Zerlegung zeigt.

b) Für $\vec{u} \in U$ und $\vec{w} \in U^\perp$ verschwindet nach Definition von U^\perp das Produkt $\vec{w} \cdot \vec{u}$, also wegen dessen (HERMITESCHER) Symmetrie auch $\vec{u} \cdot \vec{w}$. Damit ist

$$\vec{u} \in U^{\perp\perp} = \{ \vec{v} \in V \mid \vec{v} \cdot \vec{w} = 0 \text{ für alle } \vec{w} \in U^\perp \},$$

also liegt U in $U^{\perp\perp}$. Nach dem obigen Lemma ist

$$\dim U^{\perp\perp} = \dim V - \dim U^\perp = \dim V - (\dim V - \dim U) = \dim U,$$

also muß $U = U^{\perp\perp}$ sein ■

Bemerkung: Tatsächlich gilt dieses Korollar auch für unendlichdimensionale Vektorräume; da die Existenz von wie auch der Umgang mit Basen im Unendlichdimensionalen etwas problematisch ist, soll aber hier, wie bereits mehrfach in diesem Skriptum, der endlichdimensionale Fall genügen.

Definition: V sei ein endlichdimensionaler EUKLIDISCHER oder HERMITESCHER Vektorraum, und $U \leq V$ sei ein Untervektorraum. Die Abbildung $\pi_U: V \rightarrow U$, die jedem Vektor $\vec{v} = \vec{u} + \vec{w} \in V$ mit $\vec{u} \in U$ und

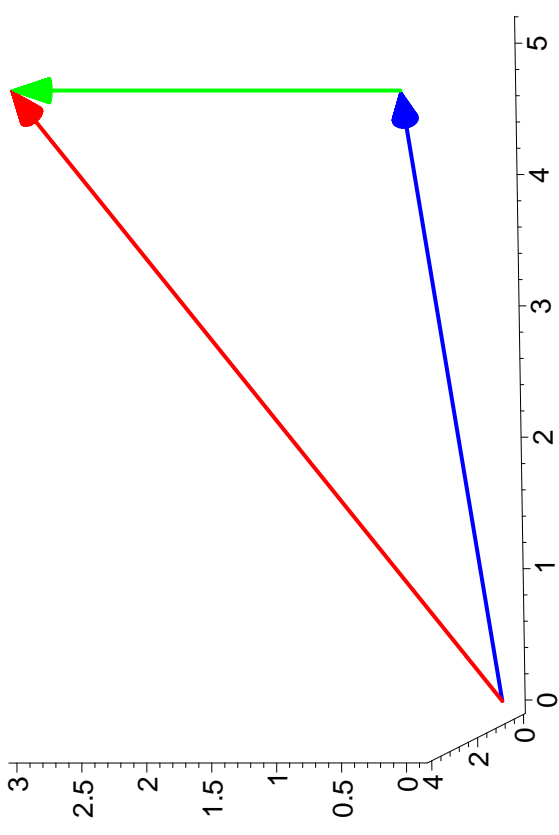


Abb. 15: Orthogonale Projektion eines Vektors

$\vec{w} \in U^\perp$ den Vektor \vec{u} zuordnet, heißt *orthogonale Projektion* von V nach U .

Wegen der eindeutigen Zerlegbarkeit eines Vektors in eine Komponente aus U und eine aus U^\perp ist π_U offensichtlich wohldefiniert und linear; der Kern von π_U ist U^\perp .

Orthogonale Projektionen sind aus der Geometrie bekannt, beispielsweise als Grundriß, Aufriß und Kreuzriß eines dreidimensionalen Körpers; uns interessiert hier vor allem ihre folgende Eigenschaft:

Lemma: V sei ein endlichdimensionaler EUKLIDISCHER oder HERMITESCHER Vektorraum, $U \leq V$ ein Untervektorraum und $\vec{v} \in V$. Dann gilt für jeden Vektor $\vec{u} \in U$ die Ungleichung $|\vec{v} - \vec{u}| \leq |\vec{v} - \pi_U(\vec{v})|$, d.h. $\pi_U(\vec{v})$ ist derjenige Vektor aus U , dessen Differenz mit \vec{v} am kürzesten ist.

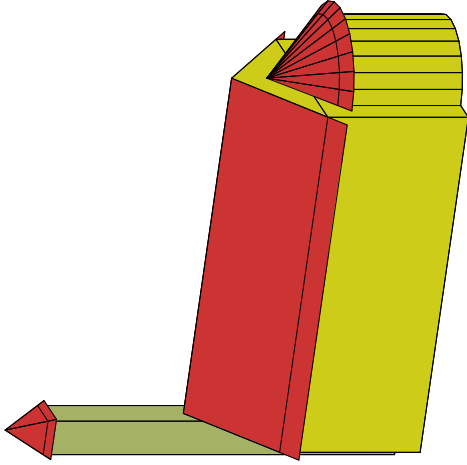


Abb. 16: Ein dreidimensionales Objekt

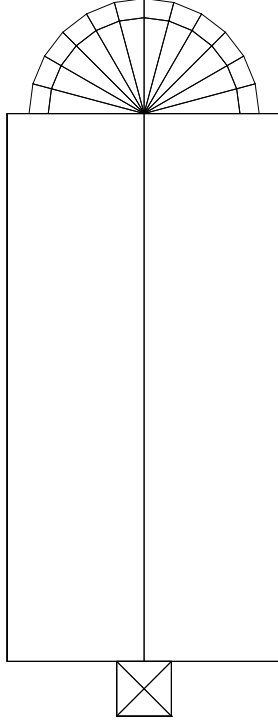


Abb. 17: Der Grundriß des Objekts aus Abbildung 16

Beweis: Wir schreiben $\vec{v} = \vec{p} + \vec{w}$ mit $\vec{p} = \pi_U(\vec{v}) \in U$ und $\vec{w} \in U^\perp$. Für jeden Vektor $\vec{u} \in U$ ist dann

$$\begin{aligned} |\vec{v} - \vec{u}|^2 &= |\vec{p} + \vec{w} - \vec{u}|^2 = |(\vec{p} - \vec{u}) + \vec{w}|^2 \\ &= (\vec{p} - \vec{u}) \cdot (\vec{p} - \vec{u}) + \vec{w} \cdot (\vec{p} - \vec{u}) + (\vec{p} - \vec{u}) \cdot \vec{w} + \vec{w} \cdot \vec{w} \\ &= |\vec{p} - \vec{u}|^2 + |\vec{w}|^2, \end{aligned}$$

denn $\vec{p} - \vec{u}$ liegt in U und \vec{w} in U^\perp . Also ist $|\vec{v} - \vec{u}|$ nie kleiner als $|\vec{v} - \vec{p}|$, und die beiden Vektoren sind genau dann gleich lang, wenn

$\vec{u} - \vec{p}$ der Nullvektor ist, also $\vec{u} = \pi_U(\vec{v})$. ■

Betrachten wir orthogonale Projektionen statt in Vektorräumen in den dazugehörigen affinen Räumen, entspricht die orthogonale Projektion auf einen Unterraum geometrisch also einfach der Konstruktion des Lotfußpunkts in diesem Unterraum.

i) Die Methode der kleinsten Quadrate

Oftmals ist zu gegebenen Beobachtungsdaten grundsätzlich bekannt, welcher Art von Gesetz sie genügen sollten; das Problem besteht „nur“ noch darin, die in diesem Gesetz vorkommenden Parameter zu bestimmen. Im einfachsten Fall könnte man etwa an einen Widerstand denken, der dadurch gemessen wird, daß man verschiedene Spannungen U_i anlegt und die zugehörigen Stromstärken I_i mißt. Nach dem Ohmschen Gesetz ist dann $U_i = R \cdot I_i$, aber aufgrund der unvermeidlichen Meßfehler werden die verschiedenen Quotienten U_i/I_i natürlich nicht alle gleich sein. Die Lösung dieses Problems ist klar: Man nimmt den Mittelwert der Quotienten. Schwieriger wird es, wenn mehrere Parameter ins Spiel kommen, wenn die Meßreihe als mehr als nur einen Parameter bestimmen soll.

Solche Fälle treten nicht nur auf in Naturwissenschaft und Technik, sondern auch in den Wirtschafts- und Sozialwissenschaften, wo es zwar selten exakte Gesetze gibt, man den Zusammenhang zwischen verschiedenen Größen aber trotzdem zumindest näherungsweise durch eine mathematische Formel beschreiben will – auch wenn diese in konkreten Einzelfällen gelegentlich ziemlich falsch sein kann.

Als Beispiel dieser Art können wir den Zusammenhang zwischen Korruption und Wohlstand in verschiedenen Staaten betrachten: edes Jahr veröffentlicht die Organisation *Transparency International* ihren *corruption perceptions index (CPI)*, in dem jedem Land eine Zahl zwischen null und zehn zugeordnet wird, je nachdem, wie stark Geschäftsleute, Risikospezialisten und die Bevölkerung die Korruption im betreffenden Land einschätzen: Ein Index von zehn bedeutet, daß es praktische keine Korruption gibt, während bei null nichts läuft ohne Bimbes. Die

neuesten Daten stammen vom 6. November 2006 und sind via

<http://www.transparency.org/>

zu finden. Die Zahlen werden als Mittelwerte über die letzten drei Jahren berechnet, so daß singuläre Ereignisse eines Jahres nicht zu sehr ins Gewicht fallen. Wir vergleichen diese Zahlen mit dem Bruttonationaleinkommen pro Einwohner, das auf dem Server des Statistischen Bundesamtes unter

http://www.destatis.de/ausl_prog/suche_ausland.htm

zu finden ist, indem man unter „Indikatoren“ das Feld „BNE je Einwohner“ auswählt. Es ist in sogenannten „Internationalen Dollar“ angegeben, das sind von der Weltbank mit einem Kaufkraftfaktor korrigierte US-\$. Die meisten Werten beziehen sich auf das Jahr 2005; für die Bahamas, Turkmenistan und die Vereinigten Arabischen Emirate sind allerdings nur ältere Daten verfügbar. In der folgenden Tabelle sind alle Staaten aufgelistet, für die sowohl das Bruttonationaleinkommen pro Einwohner als auch der CPI für 2006 vorliegt; das Bruttonationaleinkommen ist kursiv gedruckt, der Korruptionsindex fett:

Ägypten	4440	3,3
Albanien	5420	2,6
Algerien	6770	3,1
Angola	2210	2,2
Argentinien	13920	2,9
Armenien	5060	2,9
Aserbaidschan	4890	2,4
Aethiopien	1000	2,4
Australien	30610	8,7
Bahrain	21290	5,7
Bangladesch	2090	2,0
Barbados	15060	6,7
Belgien	32640	7,3
Belize	6740	3,5
Benin	1110	2,5
Bolivien	2740	2,7

Bosnien und Herzegowina	7790	2,9
Botsuana	10250	5,6
Brasilien	8230	3,3
Bulgarien	8630	4,0
Burkina Faso	1220	3,2
Burundi	640	2,4
Chile	11470	7,3
China	6600	3,3
Costa Rica	9680	4,1
Côte d’Ivoire	1490	2,1
Dänemark	33570	9,5
Deutschland	29210	8,0
Dominikanische Republik	7150	2,8
Ecuador	3070	2,3
El Salvador	5120	4,0
Eritrea	1010	2,9
Estland	15420	6,7
Finnland	31170	9,6
Frankreich	30540	7,4
Gabun	5890	3,0
Gambia	1920	2,5
Georgien	3270	2,8
Ghana	2370	3,3
Griechenland	23620	4,4
Guatemala	4410	2,6
Guinea	2240	1,9
Guyana	4230	2,5
Haiti	1840	1,8
Honduras	2900	2,5
Indien	3720	3,3
Indonesien	3460	2,4
Iran	8050	2,7
Irland	34720	7,4
Island	34760	9,6
Israel	25280	5,9
Italien	28840	4,9

Jamaika	4110	3,7	Nepal	1530	2,5
Japan	31410	7,6	Neuseeland	23030	9,6
Jemen	920	2,6	Nicaragua	3650	2,6
Jordanien	5280	5,3	Niederlande	32480	8,7
Kambodscha	2490	2,1	Niger	800	2,3
Kamerun	2150	2,3	Nigeria	1040	2,2
Kanada	32220	8,5	Norwegen	40420	8,8
Kasachstan	7730	2,6	Oman	14680	5,4
Kenia	1170	2,2	Österreich	33140	8,6
Kirgisistan	1870	2,2	Pakistan	2350	2,2
Kolumbien	7420	3,9	Panama	7310	3,1
Kongo	810	2,2	Papua-Neuguinea	2370	2,4
Kongo, Dem. Republik	720	2,0	Paraguay	4970	2,6
Korea, Republik	21850	5,1	Peru	5830	3,3
Kroatien	12750	3,4	Philippinen	5300	2,5
Kuwait	24010	4,8	Polen	13490	3,7
Laos, Dem. Volksrepublik	2020	2,6	Portugal	19730	6,6
Lesotho	3410	3,2	Ruanda	1320	2,5
Lettland	13480	4,7	Rumänien	8940	3,1
Libanon	5740	3,6	Russische Föderation	10640	2,5
Litauen	14220	4,8	Sambia	950	2,6
Luxemburg	65340	8,6	Saudi-Arabien	14740	3,3
Madagaskar	880	3,1	Schweden	31420	9,2
Malawi	650	2,7	Schweiz	37080	9,1
Malaysia	10320	5,0	Senegal	1770	3,3
Mali	1000	2,8	Sierra Leone	780	2,2
Malta	18960	6,4	Simbabwe	1940	2,4
Marokko	4360	3,2	Singapur	29780	9,4
Mauretanien	2150	3,1	Slowakei	15760	4,7
Mauritius	12450	5,1	Slowenien	22160	6,4
Mazedonien	7080	2,7	Spanien	25820	6,8
Mexiko	10030	3,3	Sri Lanka	4520	3,1
Moldau, Republik	2150	3,2	Südafrika	12120	4,6
Mongolei	2190	2,8	Sudan	2000	2,0
Mosambik	1270	2,8	Swasiland	5190	2,5
Namibia	7910	4,1	Syrien, Arabische Republik	3740	2,9

Tadschikistan	1260	2,2
Tansania, Vereinigte Republik	730	2,9
Thailand	8440	3,6
Togo	1559	2,4
Trinidad und Tobago	13170	3,2
Tschad	1470	2,0
Tschechische Republik	20140	4,8
Tunesien	7900	4,6
Türkei	8420	3,7
Turkmenistan	6910	2,2
Uganda	1500	2,7
Ukraine	6720	2,8
Ungarn	16940	5,2
Uruguay	9810	6,4
Usbekistan	2020	2,1
Venezuela	6440	2,3
Vereinigte Arabische Emirate	24090	6,2
Vereinigte Staaten	41950	7,3
Vereinigtes Königreich	32690	8,6
Vietnam	3010	2,6
Weißrussland	7890	2,1
Zentralafrikanische Republik	1140	2,4
Zypern	22230	5,6

Abbildung 18 zeigt die 147 Datenpunkte zu dieser Liste graphisch, wobei der Punkt für Deutschland etwas heller eingezeichnet ist.

Der erste Augenschein zeigt, daß korruptionsärmere Länder oftmals reicher sind: Das weitgehend korruptionsfreie Island hat ein Bruttonationaleinkommen von 34 760 \$ pro Einwohner, das deutlich korruptere Deutschland nur 29 210 \$ und ein stark korruptes Land wie Tansania nur 730 \$. Allerdings gibt es auch Ausnahmen: Beispielsweise hat Italien mit 28 840 \$ pro Einwohner zwar fast das gleiche Bruttonationaleinkommen wie Deutschland, ist aber deutlich korrupter. Es gibt also sicherlich keinen deterministischen Zusammenhang zwischen Korruption und Wohlstand, aber doch eine Tendenz.

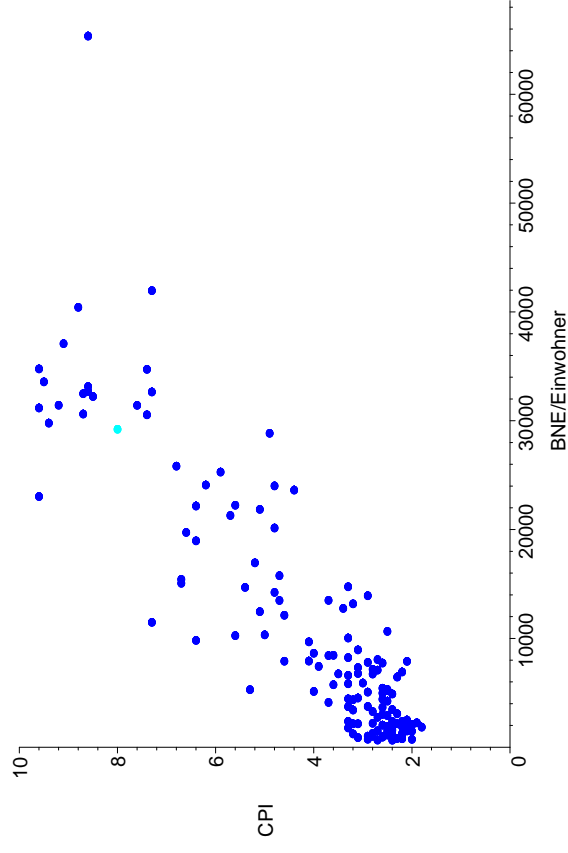


Abb. 18: Zusammenhang zwischen Korruption und Bruttonationaleinkommen je Einwohner

Falls wir nun versuchen, beispielsweise einen linearen Zusammenhang der Form

$$CPI = a + b \cdot BNE$$

zu finden, so haben wir 147 Gleichungen für die beiden unbekanntenen Koeffizienten a und b , und ein kurzer Blick auf Abbildung 18 zeigt, daß dieses lineare Gleichungssystem keine Lösung haben kann.

Wir suchen also keine Lösung, sondern zwei Zahlen a und b derart, daß die 147 Gleichungen „möglichst gut“ gelten. Was das bedeuten soll läßt sich mathematisch auf verschiedene, nicht äquivalente Weisen definieren; da wir uns im Augenblick mit Skalarprodukten beschäftigen, bietet sich an, die 147 Bruttonationaleinkommen pro Einwohner und die 147 Korruptionsindizes zu zwei Vektoren $\vec{x}, \vec{y} \in \mathbb{R}^{147}$ zusammenzufassen, und nach Zahlen a, b zu suchen, so daß die Länge des Differenzvektors $\vec{y} - a\vec{x} - b$ möglichst klein wird. Ausgeschrieben bedeutet dies, wenn wir die Komponenten von \vec{x} mit x_i und die von \vec{y}

mit y_i bezeichnen, daß die Summe

$$\sum_{i=1}^{147} (y_i - ax_i - b)^2$$

der Abweichungsquadrate möglichst klein sein soll – von daher der Name „Methode der kleinsten Quadrate“ für diesen Ansatz, mit dessen Hilfe sein Schöpfer GAUSS sowohl die Position des Planetoiden Ceres vorhersagte als auch die Vermessung und Kartierung des Königreichs Hannover durchführte.

Derselbe Ansatz läßt sich natürlich auf jedes lineare Gleichungssystem über den reellen oder komplexen Zahlen anwenden: Wir haben ein möglicherweise unlösbares lineares Gleichungssystem $A\vec{x} = \vec{b}$ und wollen einen Vektor \vec{x} so bestimmen, daß der Vektor $A\vec{x} - \vec{b}$ minimale Länge hat.

Falls das lineare Gleichungssystem lösbar ist, gibt es damit kein Problem: Wir bestimmen irgendeine Lösung \vec{x} und haben damit einen Vektor gefunden, für den $A\vec{x} - \vec{b}$ die Länge null hat – kürzer geht es nicht.

Im allgemeinen ist aber für den gesuchten Vektor \vec{x} das Produkt $A\vec{x}$ von \vec{b} verschieden; es sei etwa gleich \vec{c} . Dann ist \vec{c} ein Vektor, der sich in der Form $A\vec{x}$ darstellen läßt, und unter allen solchen Vektoren ist es derjenige, für den die Länge des Differenzvektors zu \vec{b} minimal ist. Dies erinnert an die orthogonalen Projektionen aus dem vorigen Abschnitt, und in der Tat läßt sich das Problem damit lösen:

Nehmen wir an, wir haben n Gleichungen in m Unbekannten mit Koeffizienten aus $k = \mathbb{R}$ oder $k = \mathbb{C}$. Dann definiert die Matrix $A \in k^{n \times m}$ des Gleichungssystems eine lineare Abbildung

$$\varphi: k^m \rightarrow k^n; \quad \vec{v} \mapsto A\vec{v};$$

deren Bildraum sei U . Falls die rechte Seite \vec{b} in U liegt, ist das Gleichungssystem lösbar; andernfalls suchen wir einen Vektor $\vec{x} \in k^m$, für den die Länge des Vektors $A\vec{x} - \vec{b}$ minimal wird. Da die Vektoren, die sich in der Form $A\vec{x}$ darstellen lassen, genau die Vektoren aus U sind, ist somit $A\vec{x} = \pi_U(\vec{b})$ die orthogonale Projektion von \vec{b} nach U . Diese

könnten wir im *Prinzip* bestimmen, indem wir die QR-Zerlegung von A berechnen, denn dann sind die ersten Spalten von Q eine Basis von U , die durch die weiteren Spalten zu einer Basis von ganz k^n ergänzt wird; danach haben wir ein lösbares lineares Gleichungssystem.

Wir wollen uns überlegen, wie wir \vec{x} auch ohne die rechnerisch aufwendige QR-Zerlegung bestimmen können.

Für den gesuchten Vektor \vec{x} (oder für die gesuchten Vektoren \vec{x}) ist $A\vec{x} = \varphi_U(\vec{b})$. Da $A\vec{x}$ bereits in U liegt, ist $\pi_U(A\vec{x}) = A\vec{x}$, also ist die Gleichung $A\vec{x} = \pi_U(\vec{b})$ äquivalent zu

$$\pi_U(A\vec{x}) = \pi_U(\vec{b}) \quad \text{oder} \quad A\vec{x} - \vec{b} \in \text{Kern } \pi_U = U^\perp.$$

Das orthogonale Komplement U^\perp von U besteht aus allen Vektoren $\vec{y} \in k^n$, die senkrecht stehen auf U , für die also gilt

$$(A\vec{x}) \cdot \vec{y} = 0 \quad \text{für alle } \vec{x} \in k^m.$$

Wie wir im vorletzten Abschnitt gesehen haben, ist

$$(A\vec{x}) \cdot \vec{y} = \vec{x} \cdot A^* \vec{y} \quad \text{für alle } \vec{x} \in k^m, \vec{y} \in k^n,$$

\vec{y} liegt also genau dann in U^\perp , wenn $A^* \vec{y}$ senkrecht steht auf allen Vektoren $\vec{x} \in k^m$. Ein solcher Vektor aus k^m ist insbesondere $A^* \vec{y}$ selbst; wegen der positiven Definitheit des (HERMITESCHEN) Skalarprodukts ist also $A^* \vec{y} = \vec{0}$. Da aus $A^* \vec{y} = \vec{0}$ für alle $\vec{x} \in k^m$ folgt, daß $\vec{x} \cdot A^* \vec{y}$ verschwindet, ist damit

$$U^\perp = \{ \vec{y} \in k^n \mid A^* \vec{y} = \vec{0} \}.$$

$A\vec{x} - \vec{b}$ liegt also genau dann im Kern von π_U , wenn $A^*(A\vec{x} - \vec{b}) = \vec{0}$ ist oder, anders ausgedrückt, wenn \vec{x} eine Lösung des linearen Gleichungssystems

$$(A^* A) \vec{x} = A^* \vec{b}$$

ist. Da die adjungierte Matrix A^* einfach die transponierte Matrix zur komplex konjugierten Matrix zu A ist, wobei die komplexe Konjugation über \mathbb{R} natürlich entfällt, läßt sich dieses Gleichungssystem schnell aufstellen und dann nach GAUSS lösen.

Betrachten wir dies konkret im eingangs diskutierten Fall eines linearen Zusammenhangs $y = ax + b$ zu N Wertepaaren $(x_i, y_i) \in \mathbb{R}^2$, wobei N sinnvollerweise größer als zwei sein sollte. Wir haben dann N Gleichungen

$$y_i = ax_i + b \quad \text{oder} \quad x_i a + b = y_i,$$

wobei hier im Gegensatz zu unserer sonstigen Gewohnheit die Parameter a und b unbekannt sind, während die x_i und die y_i bekannt sind. Wir haben also ein lineares Gleichungssystem von N Gleichungen in den beiden Variablen a und b .

Fassen wir die Werte x_i zusammen zu einem Vektor $\vec{x} \in \mathbb{R}^N$ und die y_i zu einem Vektor $\vec{y} \in \mathbb{R}^n$, so läßt sich dieses Gleichungssystem kurz schreiben als

$$\vec{x} \cdot a + \vec{1} \cdot b = \vec{y},$$

wobei $\vec{1} \in \mathbb{R}^N$ jenen Vektor bezeichnen soll, dessen sämtliche Komponenten eins sind.

Die Matrix des Gleichungssystems ist somit die $N \times 2$ -Matrix A mit Spalten \vec{x} und $\vec{1}$. Da wir mit reellen Zahlen rechnen, ist A^* einfach die transponierte Matrix dazu, also die $2 \times N$ -Matrix, in deren erster Zeile die x_i stehen, während in der zweiten lauter Einsen stehen. Somit ist

$${}^t A A = \begin{pmatrix} \vec{x} \cdot \vec{x} & \vec{x} \cdot \vec{1} \\ \vec{x} \cdot \vec{1} & \vec{1} \cdot \vec{1} \end{pmatrix} \quad \text{und} \quad {}^t A \vec{b} = \begin{pmatrix} \vec{x} \cdot \vec{y} \\ \vec{1} \cdot \vec{y} \end{pmatrix},$$

das Gleichungssystem wird also zu

$$(\vec{x} \cdot \vec{x})a + (\vec{x} \cdot \vec{1})b = \vec{x} \cdot \vec{y} \quad \text{und} \quad (\vec{x} \cdot \vec{1})a + N b = \vec{1} \cdot \vec{y}.$$

Seine Matrix ist genau dann singulär, wenn die Determinante verschwindet, wenn also $N(\vec{x} \cdot \vec{x}) = (\vec{x} \cdot \vec{1})^2$ ist. Nach der CAUCHY-SCHWARZschen Ungleichung ist

$$|\vec{1} \cdot \vec{x}| \leq |\vec{1}| \cdot |\vec{x}| = \sqrt{N} |\vec{x}|, \quad \text{also} \quad |\vec{1} \cdot \vec{x}|^2 \leq N(\vec{x} \cdot \vec{x})$$

mit Gleichheit nur dann, wenn die Vektoren \vec{x} und $\vec{1}$ linear abhängig sind, wenn also alle x_i denselben Wert x haben. In diesem Fall ist die erste Gleichung das x -fache der zweiten, es gibt also unendlich viele Lösungen.

Andernfalls ist die Matrix invertierbar, die Lösung also eindeutig.

Führen wir die (in der Ausgleichsrechnung ziemlich verbreiteten) Abkürzungen

$$[x^T] = \sum_{i=1}^N x_i^T, \quad [y^T] = \sum_{i=1}^N x_i^T y_i^T \quad \text{und} \quad [x^T y^T] = \sum_{i=1}^N x_i^T y_i^T$$

ein, so erhält das Gleichungssystem die übersichtlichere Gestalt

$$[x^T]a + [x^T]b = [xy] \quad \text{und} \quad [x]a + Nb = [y].$$

Subtraktion von $[x]/[x^2]$ mal der ersten Gleichung von der zweiten führt auf

$$\left(N - \frac{[x]^2}{[x^2]} \right) b = [y] - \frac{[x]}{[x^2]} [xy]$$

oder $(N[x^2] - [x]^2)b = [y][x^2] - [x][xy]$, d.h.

$$b = \frac{[y][x^2] - [x][xy]}{N[x^2] - [x]^2}.$$

(Man beachte, daß im Falle der eindeutigen Lösbarkeit sowohl $[x^2] > 0$ als auch $N[x^2] - [x]^2 > 0$ ist.)

Einsetzen von b in die erste Gleichung ergibt dann auch

$$a = \frac{[xy] - [x]b}{[x^2]}.$$

Im Falle des Zusammenhangs zwischen Korruptionsindex CPI und Bruttonationaleinkommen pro Einwohner BNE erhalten wir nach diesen Formeln die Ausgleichsgerade

$$\text{CPI} = 2,29265 + 0,00017682 \cdot \text{BNE},$$

die Steigung ist also erwartungsgemäß positiv. Der relativ große konstante Term zeigt, daß *im Mittel* Korruption selbst bei sehr armen Ländern deutlich über dem unteren Ende der Skala liegt. Abbildung 19 zeigt die Ausgleichsgerade zusammen mit den Daten.

Natürlich sind die Datenpunkte relativ breit gestreut um die Ausgleichsgerade; der Zusammenhang zwischen Korruption und Wohlstand ist

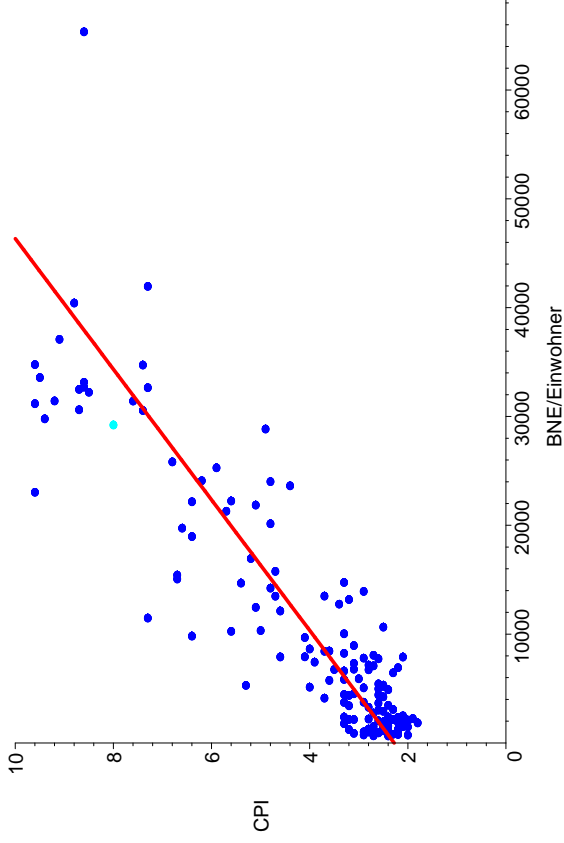


Abb. 19: Ausgleichsgerade zu Abbildung 18

schließlich zum Glück kein unausweichliches deterministisches Gesetz, sondern nur eine empirische Beobachtung.

Auch bei Messungen physikalischer Größen, wo die verschiedenen Meßgrößen meist durch wohlbekannte Naturgesetze miteinander verbunden sind, gibt es praktisch immer eine Streuung der Daten um die theoretisch richtige Meßkurve; absolut fehlerfreie Messungen sind, trotz aller Mühe der Experimentatoren, fast nie möglich, da es praktisch immer ein Grundrauschen der Meßgeräte und/oder nicht in ihrer Gesamtheit erfassbare Umgebungseinflüsse u_{sw} gibt. Vor allem bei Messungen, mit denen Konstanten für Naturgesetze ermittelt werden sollen oder gar ein Experiment zwischen zwei oder mehr Hypothesen entscheiden soll, ist es daher wichtig zu wissen, wie gut die Übereinstimmung zwischen den Daten und der berechneten Kurve (oder Fläche u_{sw} .) wirklich ist.

Solche Maße stellt die Statistik zur Verfügung; für ihr Verständnis sind daher meist zumindest Grundlagenkenntnisse der Statistik notwendig, wie wir sie (wenn auch nur kurz) im nächsten Semester behandeln

werden. Im einfachsten und zugleich wichtigsten Fall eines linearen Zusammenhangs zwischen zwei Größen allerdings reicht die lineare Algebra, um das sowohl in der Theorie wie auch den Anwendungen wichtigste Qualitätsmaß zu definieren, den Korrelationskoeffizienten.

Angenommen, wir haben N Datenpaare (x_i, y_i) , zwischen denen ein perfekter linearer Zusammenhang besteht, d.h.

$$y_i = ax_i + b \quad \text{für alle } i = 1, \dots, N.$$

Wir wollen den Datenvektoren $\vec{x} \in \mathbb{R}^N$ mit Komponenten x_i und $\vec{y} \in \mathbb{R}^N$ mit Komponenten y_i Vektoren zuordnen, die nicht nur in einem linearen Zusammenhang stehen, sondern sogar gleich sind; mit anderen Worten, wir wollen die Parameter a und b aus obiger Gleichung eliminieren.

Dazu betrachten wir als erstes die Mittelwerte

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i \quad \text{und} \quad \bar{y} = \frac{1}{N} \sum_{i=1}^N y_i.$$

Da $y_i = ax_i + b$ ist für alle i , folgt sofort, daß auch $\bar{y} = a\bar{x} + b$ ist, und damit $(y_i - \bar{y}) = a(x_i - \bar{x})$ für alle $i = 1, \dots, N$.

Damit ist der Parameter b eliminiert. Bezeichnen wir wieder mit $\vec{1} \in \mathbb{R}^N$ den Vektor, dessen sämtliche Komponenten Einsen sind, ist nun also $(\vec{y} - \bar{y}\vec{1}) = a(\vec{x} - \bar{x}\vec{1})$. Aus dieser Gleichung können wir nun leicht a bis auf sein Vorzeichen eliminieren, indem wir die beiden Vektoren durch ihre Länge dividieren. Dies ist natürlich nur möglich, wenn keiner der beiden Vektoren gleich dem Nullvektor ist, wenn also nicht alle $x_i = \bar{x}$ oder alle $y_i = \bar{y}$ sind. Bei nicht getürkten Messungen ist dies allerdings *praktisch* nie der Fall, so daß die Nützlichkeit der folgenden Diskussion und Definition nicht darunter leidet, daß wir diesen Fall ausschließen müssen.

Falls also weder $\vec{y} - \bar{y}\vec{1}$ noch $\vec{x} - \bar{x}\vec{1}$ der Nullvektor ist, betrachten wir die beiden auf Länge eins normierten Vektoren

$$\frac{\vec{y} - \bar{y}\vec{1}}{\|\vec{y} - \bar{y}\vec{1}\|} \quad \text{und} \quad \frac{\vec{x} - \bar{x}\vec{1}}{\|\vec{x} - \bar{x}\vec{1}\|}.$$

Diese sind nun offensichtlich entweder gleich (für $a > 0$) oder entgegengesetzt gleich (für $a < 0$).

Wenn (wie in der Realität meist der Fall) *kein* perfekter linearer Zusammenhang zwischen den x_i und den y_i besteht, können wir trotzdem – falls weder $\vec{y} - \bar{y}\vec{1}$ noch $\vec{x} - \bar{x}\vec{1}$ der Nullvektor ist – die beiden Vektoren

$$\frac{\vec{y} - \bar{y}\vec{1}}{|\vec{y} - \bar{y}\vec{1}|} \quad \text{und} \quad \frac{\vec{x} - \bar{x}\vec{1}}{|\vec{x} - \bar{x}\vec{1}|}$$

betrachten. Da beides Einheitsvektoren sind, unterscheiden sie sich nur in der Richtung; als Maß für ihren Unterschied bietet sich daher den Winkel zwischen \vec{x} und \vec{y} an. Rechnerisch einfacher ist der Cosinus dieses Winkels, denn der ist bei Einheitsvektoren einfach gleich dem Skalarprodukt.

Definition: Der Korrelationskoeffizient zwischen zwei Datenvektoren \vec{x} und $\vec{y} \in \mathbb{R}^n$, die keine Vielfachen des Vektors $\vec{1} \in \mathbb{R}^n$ sind, ist

$$\rho = \frac{(\vec{x} - \bar{x} \cdot \vec{1}) \cdot (\vec{y} - \bar{y} \cdot \vec{1})}{|\vec{x} - \bar{x} \cdot \vec{1}| \cdot |\vec{y} - \bar{y} \cdot \vec{1}|}.$$

Damit ist also $\rho = \pm 1$ genau dann, wenn es einen perfekten linearen Zusammenhang $y_i = ax_i + b$ zwischen den beiden Größen gibt, mit $\rho = 1$ für $a > 0$ und $\rho = -1$ für $a < 0$. Ansonsten ist der Zusammenhang umso besser, je größer der Betrag von ρ ist. Für $\rho = 0$ stehen die beiden Vektoren $\vec{x} - \bar{x}\vec{1}$ und $\vec{y} - \bar{y}\vec{1}$ senkrecht aufeinander, d.h. wenn x_i größer ist als der Mittelwert \bar{x} , kann y_i im Mittel genauso gut größer wie auch kleiner als der Mittelwert \bar{y} sein. (in der Statistik ist dies die *Definition* für die Unabhängigkeit von Daten.)

Definition: Zwei Größen x und y heißen $\begin{cases} \text{positiv} \\ \text{negativ} \end{cases}$ korreliert, wenn $\rho \begin{cases} \geq \\ < \end{cases} 0$ ist. Sie heißen unkorreliert oder voneinander unabhängig, wenn $\rho = 0$ ist.

Im Beispiel der Korruption erhalten wir einen Korrelationskoeffizienten von $\rho \approx 0,885395$; dies entspricht einem Winkel von etwa $27,7^\circ$ zwischen den oben definierten Vektoren.

Um ein Gefühl für Korrelationskoeffizienten zu bekommen, wollen wir zwei Beispiele betrachten, die sich zumindest visuell sehr unterscheiden: Der CPI für Deutschland hatte in den letzten Jahren folgende Werte:

Jahr:	1980–1985	1988–1992	1995	1996	1997	1998	1999
CPI:	8,14	8,13	8,14	8,27	8,23	7,9	8,0
Jahr:	2000	2001	2002	2003	2004	2005	2006
CPI:	7,6	7,4	7,3	7,7	8,2	8,2	8,0

Wie Abbildung 20 zeigt, sieht der Zusammenhang zwischen Jahr und CPI nicht sonderlich linear aus: Der Bimbesknick ist unverkennbar, jedoch scheint die Talsohle inzwischen durchschritten, so daß die abwärts gehende Ausgleichsgerade trotz des erneuten Abfalls hoffentlich nicht den derzeitigen Trend beschreibt. Der Korrelationskoeffizient $\kappa \approx -0,317$ ist erwartungsgemäß ziemlich schlecht: Er ist der Kosinus eines Winkels von etwa $108,5^\circ$.

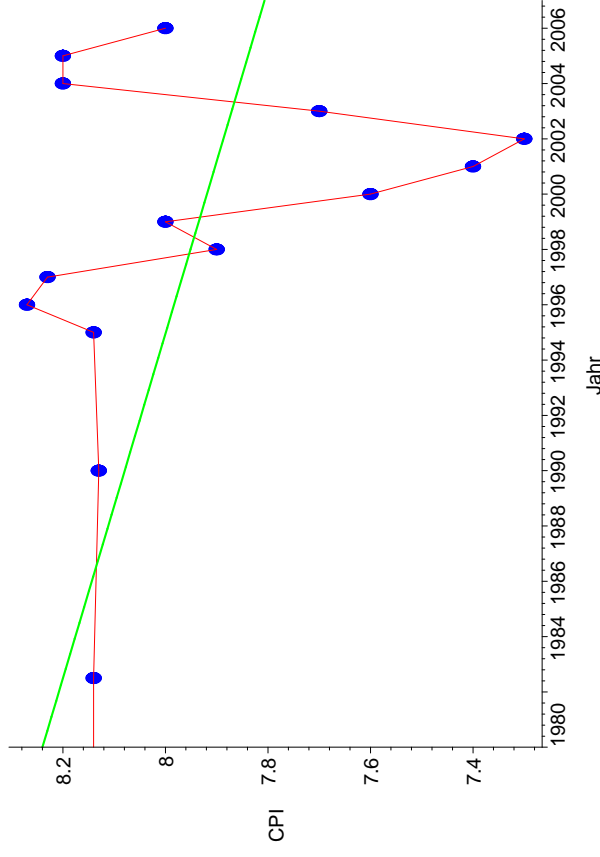


Abb. 20: Zeitabhängigkeit des CPI für Deutschland

Vergleichen wir dagegen die Mannheimer Ergebnisse von Europawahl und Gemeinderatswahl vom 13. Juni 2004 miteinander, so gibt es bei keiner der vier Parteien, die zu beiden Wahlen angetreten ist, dramatische Unterschiede zwischen ihrem Stimmanteil bei den beiden Wahlen, obwohl gewisse Abweichungen unverkennbar sind.

	Europawahl	Gemeinderatswahl
CDU	38,14	40,41
SPD	28,91	33,38
Grüne	14,72	10,19
FDP	5,86	3,43

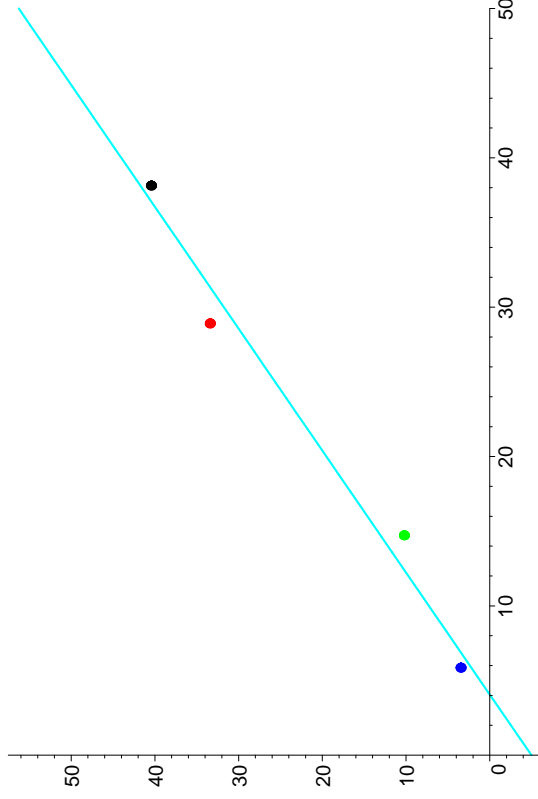


Abb. 21: Zusammenhang Europawahl/Gemeinderatswahl

Wie Abbildung 21 zeigt, kann man den Zusammenhang in recht guter Näherung durch eine Gerade beschreiben, und in der Tat erhalten wir hier $\kappa \approx 0,9900614$, was einem Winkel von etwa acht Grad entspricht.

Korrelationskoeffizienten mit kleinem Betrag müssen nicht unbedingt bedeuten, daß kein deterministischer Zusammenhang zwischen den Daten besteht: Sie besagen nur, daß es keinen *linearen* Zusammenhang gibt. Betrachtet man etwa Wertepaare $(x_i, \sin x_i)$, so erhält man für Werte x_i , die einigermaßen gleichmäßig über eine oder mehrere Perioden der Sinusfunktion verteilt sind, einen Korrelationskoeffizienten nahe Null, obwohl der Zusammenhang zwischen den beiden werten eines jeden Paares strikt deterministisch ist. In so einem Fall ist einfach der lineare Ansatz die falsche Strategie und man muß alternative Ansätze finden.

j) Euklidische Vektorräume in der Informationssuche

In §4j) haben wir gesehen, wie Google die Unzahl von Internetseiten völlig unabhängig von jeder Suchanfrage nach ihrer Wichtigkeit ordnet. Hier soll es kurz skizziert werden, wie die Lineare Algebra auch hilft beim Problem, zu einer Suchanfrage geeignete Dokumente zu finden. Dabei muß es nicht unbedingt um Suche im Internet gehen; mindestens genauso wichtig sind wissenschaftliche Literaturdatenbanken, in denen zumindest Tausende (meist deutlich mehr) wissenschaftlicher Arbeiten gespeichert sind, aus denen ein Anwender die für seine Forschung relevanten finden möchte. Am einfachsten geht das, wenn entweder der Autor oder ein Berichterstatter die Arbeiten nach Themengebieten ordnet: In der Mathematik etwa gibt es dazu ein umfangreiches Klassifikationsschema der beiden westlichen Referatorgane *Zentralblatt für Mathematik* und *ihre Grenzgebiete* und *Mathematical Reviews*, das auch die meisten Fachzeitschriften verwenden; in anderen Wissenschaften ist es ähnlich.

Solche Zuordnungen sind meistens recht genau, da sie von Experten der jeweiligen *Teilgebiete* vorgenommen werden; andererseits birgt natürlich auch gerade das die Gefahr in sich, daß eine Arbeit, die für mehrere Gebiete relevant ist, möglicherweise nur denen zugeordnet wird, für die sich der Autor oder Berichterstatter interessiert. Außerdem ist selbst eine sehr detaillierte Einteilung, die im Falle der Mathematik immerhin 35 Seiten Kleingedrucktes benötigt, immer noch zu grob, um genau *die* drei Arbeiten zu finden, in denen ein sehr spezielles Problem behandelt wird. Im Internet mit seiner Vielzahl von teilweise sehr schnell vari-

ierenden Informationsangeboten ist ein solcher Ansatz von vornherein chancenlos.

Zusätzlich zur Klassifikation durch menschliche Experten braucht man daher bei der Informationssuche auch Algorithmen, die gelesene Informationen automatisch klassifizieren und bezüglich ihrer Relevanz zu einer konkreten Suchanfrage beurteilen können.

Große Internetsuchmaschinen verwenden dazu eine Vielzahl von Algorithmen; mit Ausnahme von [google.com](http://www.google.com), die ihr Rangbildungsverfahren unter

<http://www.google.com/technology/pigeonrank.html>

mehr oder weniger ausführlich beschreiben, schweigen sie sich allerdings aus über die genauen Einzelheiten und Parameter: Schließlich sollen die vielen unseriösen Anbieter, die mit allen Tricks Besucher auf ihre Webseiten locken wollen, nicht auch noch unterstützt werden.

Wir müssen uns daher auf die grundlegenden mathematischen Algorithmen beschränken, die wohl in der einen oder anderen Form in praktisch jeder Suchmaschine zu finden sind und die, als Gegenstand wissenschaftlicher Forschung, natürlich öffentlich bekannt sind.

Die ersten Systeme arbeiteten mit den üblichen Suchalgorithmen aus der Textverarbeitung, durchsuchten also alle gespeicherten Dokumente nach dem Vorkommen einer oder mehrerer vorgegebener Zeichenketten. Auch wenn es dafür sehr effiziente Algorithmen gibt, ist dieses Verfahren bei wirklich großen Datenmengen nicht mehr mit realistischem Aufwand durchführbar, so daß nun meist Verfahren aus der linearen Algebra verwendet werden.

Dazu wird eine Liste von Suchbegriffen s_i , $i = 1, \dots, n$ festgelegt – beispielsweise die Wörter aus einem Wörterbuch der Dokumentenpraxis. Oftmals werden darauf noch geeignete Operationen angewandt wie *stemming*, d.h. Wörter mit gleichem Stamm werden miteinander identifiziert, oder *latent semantic indexing*, wo durch Clusterbildung bei den vorhandenen Dokumenten Begriffspaare identifiziert werden, die im allgemeinen im gleichen Kontext auftreten und die dann auch bei

Suchanfragen als äquivalent betrachtet werden; außerdem werden sogenannte „Nullwörter“, die für Suchanfragen typischerweise ohne Bedeutung sind, eliminiert. Dabei handelt es sich beispielsweise um Artikel und Praepositionen, gelegentlich aber auch um spezifische Wörter aus dem Kontext des jeweiligen Systems: Bei Boeing, die ein solches System zur Verwaltung ihrer Wartungshandbücher aufbauten, ist etwa das Wort „aeroplane“ ein Nullwort – die Gesellschaft verkauft schließlich keine Rasenmäher.

Sind nun m Dokumente zu betrachten, so bildet man eine $n \times m$ -Matrix A , deren Eintrag a_{ij} etwas über das Vorkommen des i -ten Suchbegriffs im j -ten Dokument aussagt. Im einfachsten Fall setzt man einfach $a_{ij} = 1$, falls der Begriff vorkommt und null sonst, alternativ kann a_{ij} auch die Häufigkeit des Begriff im Dokument sein, wobei diese Häufigkeit oft noch gewichtet wird, indem beispielsweise Vorkommen im vorderen Teil des Dokuments höher gewichtet wird oder aber die Suchmaschine ohnehin nur den Anfangsteil des Dokuments bis zu einer gewissen Maximallänge berücksichtigt. Auch das Vorkommen in Überschriften oder zwischen $\langle \text{META} \rangle$ -tags kann eventuell gesondert behandelt werden, indem man beispielsweise Inhalte, die im Browserfenster nicht sichtbar werden, wegen der damit verbundenen Mißbrauchsmöglichkeit ignoriert. Gelegentlich wird auf das Ergebnis noch eine Skalierungsfunktion wie etwa $\log(1 + x)$ angewendet.

Die entstehende Matrix ist natürlich riesig; schon 1998 wurde geschätzt, daß allein für englischsprachige Dokumente bis zu 300 000 Suchbegriffe notwendig sind, die in etwa 300 Millionen Dokumenten gesucht werden müssen; die Matrix hat also knapp hundert Billionen Einträge. Bei nur einem Byte pro Eintrag hätte man also bei der Speicherung als Feld einen Platzbedarf von etwa 90 Terabyte.

Nun kommt allerdings in fast jedem Dokument nur ein verschwindend geringer Bruchteil der Suchbegriffe vor, so daß die meisten Einträge von A Nullen sind. Die Matrix läßt sich daher erheblich kompakter speichern, wenn man beispielsweise nur die Tripel (i, j, a_{ij}) notiert, für die $a_{ij} \neq 0$ ist. Die numerische Mathematik kennt eine ganze Reihe von Algorithmen, mit denen man auch solche sogenannte „spärlich besetzte“ Matrizen effizient behandeln kann.

Der Inhalt des j -ten Dokuments wird nun also kodiert durch den j -ten Spaltenvektor der Matrix A , einen Vektor aus \mathbb{R}^n . Auch eine Suchanfrage läßt sich durch einen solchen Vektor kodieren, indem man die j -te Komponente auf eins setzt, falls der j -te Suchbegriff in der Anfrage vorkommt, und auf null sonst. (Man kann natürlich auch andere Werte wählen und beispielsweise seltene Wörter höher gewichten als häufige LSW.)

Ein Dokument sollte umso besser zu einer Suchanfrage passen, je weniger sich die dazu gehörigen Vektoren voneinander unterscheiden. Als Maß für den Unterschied zweier Vektoren haben wir im vorigen Abschnitt den Cosinus des eingeschlossenen Winkels kennengelernt; falls man die Spaltenvektoren der Matrix auf Länge Eins normiert, läßt sich dieser durch eine einziges Skalarprodukt berechnen. Ein Dokument wird dann als relevant für die Suchanfrage betrachtet, wenn dieser Wert über einer festzulegenden Schranke liegt, und die so gefundenen Dateien können dann eventuell noch mit anderen Methoden (Volltextsuche, Links von anderen Seiten, ...) weiter untersucht werden zur Festlegung der endgültigen Reihenfolge, in der sie dem Benutzer gezeigt werden.

Für sehr große Datenmengen ist allerdings die Matrix A trotz ihrer spärlichen Besetzung immer noch zu groß; wie bei der Komprimierung von Bilddaten sucht man daher nach einer Art und Weise, sie bei möglichst geringem Informationsverlust deutlich zu komprimieren. Ein angenehmer Nebeneffekt dabei ist, wie experimentelle Untersuchungen zeigen, auch eine gewisse „Rauschunterdrückung“: Es ist zwar schwierig, exakt zu definieren, was „Rauschen“ in einer Term-Dokument-Matrix sein soll, aber jeder wird wohl damit übereinstimmen, daß etwa dieses Skriptum nicht die ideale Referenz zum Thema „Rasenmäher“ ist, obwohl dieses Wort hier nun schon zum zweiten Mal vorkommt.

Einen Ansatz zur Datenreduktion liefert die QR -Zerlegung: Ist $A = QR$ und $\vec{a} \in \mathbb{R}^n$ eine Suchanfrage, so ist für die j -ten Spalten \vec{a}_j von A und \vec{r}_j von R

$$\vec{a} \cdot \vec{a}_j = \vec{a} \cdot (Q\vec{r}_j) = (Q\vec{a}) \cdot \vec{r}_j,$$

und die Matrix R wird im allgemeinen deutlich mehr Nullen enthalten

als A , da der Rang von A wohl deutlich unter n liegen dürfte. Eine weitere Komprimierung wird dadurch erreicht, daß Einträge von R , die unterhalb einer gewissen Schranke liegen, auf Null gesetzt werden; dadurch ändert sich bei hinreichend kleiner Schranke an den meisten Skalarprodukten nicht viel, dafür verringert sich aber der Speicherbedarf noch einmal beträchtlich.

Oft verwendet man anstelle der QR -Zerlegung auch die hier nicht behandelte Singulärwertzerlegung von A : Danach läßt sich A schreiben als Produkt UDV mit orthogonalen Matrizen $U \in \mathbb{R}^{n \times n}$ und $V \in \mathbb{R}^{m \times m}$ sowie einer Diagonalmatrix $D \in \mathbb{R}^{n \times m}$. U und V können so gewählt werden, daß die Diagonaleinträge der Größe nach angeordnet sind, und man erhält die gewünschte Rangreduktion, indem man alle Einträge unterhalb einer gewissen Größe auf Null setzt.

Eine ausführlichere Darstellung der Verfahren zur Textsuche, die keine über den Inhalt dieses Skriptums hinausgehende Mathematikkenntnisse voraussetzt, findet man beispielsweise in

MICHAEL W. BERRY, MURRAY BROWNE: Understanding Search Engines: Mathematical Modeling and Text Retrieval, SIAM, 1999, ²2005,

Fallstudien im Tagungsband

MICHAEL W. BERRY [Hrsg.]: Computational Information Retrieval, SIAM, 2001.