

# Einführung in die Hauptkomponentenanalyse

Florian Steinke

16. Juni 2009

## 1 Vorbereitung: Einige Aspekte der multivariaten Gaußverteilung

**Definition 1.1.** Die 1-D Gaußverteilung für  $x \in \mathbb{R}$  ist

$$p(x) \propto \exp\left(-\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2}\right). \quad (1)$$

Notation:  $x \sim N(x; \mu, \sigma^2)$ .

Eigenschaften:

- Mittelwert  $E(x) = \mu$ , Varianz  $E((x - E(x))^2) = \sigma^2$ .

**Beispiel 1.2.** Sei  $\mathbf{x} \in \mathbb{R}^2$  und

$$\begin{aligned} p(\mathbf{x}) &= N(x_1; \mu_1, \sigma_1^2) N(x_2; \mu_2, \sigma_2^2) \\ &\propto \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \begin{pmatrix} \sigma_1^{-2} & \\ & \sigma_2^{-2} \end{pmatrix} (\mathbf{x} - \boldsymbol{\mu})\right) \end{aligned}$$

**Beispiel 1.3.** Sei  $\mathbf{x} \in \mathbb{R}^2$ ,  $\mathbf{U} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix}$ , und

$$\begin{aligned} p(\mathbf{x}) &\propto \exp\left(-\frac{1}{2} (\mathbf{U}(\mathbf{x} - \boldsymbol{\mu}))^T \begin{pmatrix} \sigma_1^{-2} & \\ & \sigma_1^{-2} \end{pmatrix} \mathbf{U}(\mathbf{x} - \boldsymbol{\mu})\right) \\ &\propto \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \underbrace{\mathbf{U}^T \begin{pmatrix} \sigma_1^{-2} & \\ & \sigma_1^{-2} \end{pmatrix} \mathbf{U}}_{\equiv \boldsymbol{\Sigma}^{-1}} (\mathbf{x} - \boldsymbol{\mu})\right) \end{aligned}$$

Also  $\boldsymbol{\Sigma} = \mathbf{U} \begin{pmatrix} \sigma_1^2 & \\ & \sigma_1^2 \end{pmatrix} \mathbf{U}^T$ .

**Definition 1.4.** Die M-D Gaußverteilung für  $\mathbf{x} \in \mathbb{R}^M$  ist

$$p(\mathbf{x}) \propto \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right). \quad (2)$$

Notation:  $\mathbf{x} \sim N(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ .

Eigenschaften:

- Mittelwert

$$E(\mathbf{x}) = \boldsymbol{\mu} \quad (3)$$

Kovarianz

$$E((\mathbf{x} - E(\mathbf{x}))(\mathbf{x} - E(\mathbf{x}))^T) = \boldsymbol{\Sigma}. \quad (4)$$

- $\Sigma$  muss symmetrisch positiv definit (spd) sein.

*Proof.* Symmetrisch ist klar nach (4). Zum Beweis, dass  $\Sigma$  positiv definit ist, wähle  $\mathbf{a} \neq 0 \in \mathbb{R}^M$ . Sei o.B.d.A.  $\boldsymbol{\mu} = 0$ . Dann ist nach (4)

$$\mathbf{a}^T \Sigma \mathbf{a} = E(\mathbf{a}^T \mathbf{x} \mathbf{x}^T \mathbf{a}) = E((\mathbf{a}^T \mathbf{x})^2) \geq 0.$$

Falls  $E((\mathbf{a}^T \mathbf{x})^2) = 0$ , dann ist mit Wahrscheinlichkeit eins  $\mathbf{x} = 0$  in Richtung  $\mathbf{a}$ , d.h. der Ellipsoid der Gaußverteilung ist **entartet** in Richtung  $\mathbf{a}$ .  $\square$

- Jede Gaußverteilung läßt sich wie in Beispiel 1.3 darstellen.

*Proof.* Da  $\Sigma$  spd, hat es eine Eigenwertdarstellung mit orthogonalen Matrizen, d.h. man kann schreiben  $\Sigma = \mathbf{U} \mathbf{D} \mathbf{U}^T$ , wobei  $\mathbf{U}$  orthogonal, d.h.  $\mathbf{U} \mathbf{U}^T = \mathbf{U}^T \mathbf{U} = \mathbf{1}$ , und  $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_M)$  diagonal. Da  $\Sigma$  positiv definit, sind die Eigenwerte  $\lambda_i > 0$  und  $\sigma_i^2 = \lambda_i$ .  $\square$

Bemerkung:  $\mathbf{U}, \mathbf{D}$  sind bis auf Permutationen eindeutig bestimmt (falls alle Eigenwerte verschieden). Wir wählen typischerweise die Darstellung  $\lambda_1 > \lambda_2 > \dots > \lambda_M$ .

- Alle 1-D Schnitte einer M-D Gaußverteilung sind 1-D Gaußverteilt, d.h. für  $\mathbf{x} = \mathbf{a} + t\mathbf{b}$  ist  $t$  immer 1-D Gaußverteilt.

*Proof.* (für den Fall  $\mathbf{a} = \boldsymbol{\mu} = 0$ ). Es ist

$$\begin{aligned} p(\mathbf{x}) d\mathbf{x} &= p(\mathbf{x}_t) \frac{d\mathbf{x}}{dt} dt \\ &\propto \exp\left(-\frac{1}{2} \mathbf{b}^T \Sigma^{-1} \mathbf{b} t^2\right) dt = N(t; 0, (\mathbf{b}^T \Sigma^{-1} \mathbf{b})^{-1}) dt. \end{aligned}$$

$\square$

- Für lineare Funktionen  $\mathbf{y} = \mathbf{A}\mathbf{x}$  gilt: Aus  $\mathbf{x} \sim N(\mathbf{x}; \boldsymbol{\mu}, \Sigma)$  folgt  $\mathbf{y} \sim N(\mathbf{y}; \mathbf{A}\boldsymbol{\mu}, \mathbf{A}\Sigma\mathbf{A}^T)$ .
- Mit Rauschen  $\boldsymbol{\epsilon} \sim N(\boldsymbol{\epsilon}; 0, \Sigma')$  und  $\mathbf{x} \sim N(\mathbf{x}; \boldsymbol{\mu}, \Sigma)$  gilt für  $\mathbf{y} = \mathbf{x} + \boldsymbol{\epsilon}$ , dass  $\mathbf{y} \sim N(\mathbf{y}; \boldsymbol{\mu}, \Sigma + \Sigma')$ .

Kleiner Ausblick auf weitere wichtige Eigenschaften der Gaußverteilung:

- Konditionierung und Marginalisierung: Sei  $\mathbf{x} = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix}$ ,  $\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}$ ,  $\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$  und  $\mathbf{x} \sim N(\mathbf{x}; \boldsymbol{\mu}, \Sigma)$ . Dann ist

$$p(\mathbf{x}_1) = \int p(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_2 = N(\mathbf{x}_1; \boldsymbol{\mu}_1, \Sigma_{11}) \quad (5)$$

$$p(\mathbf{x}_1 | \mathbf{x}_2) = N(\mathbf{x}_1; \boldsymbol{\mu}_1 + \Sigma_{12} \Sigma_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2), \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}) \quad (6)$$

- Die Familie der Gaußverteilungen ist geschlossen unter Multiplikation (Division), Faltung, Marginalisierung und Conditionierung.
- $N(\mathbf{x}; \boldsymbol{\mu}, \Sigma)$  ist die Verteilung, die gegeben den Mittelwert  $\boldsymbol{\mu}$  und die Kovarianz  $\Sigma$ , die maximale Entropie hat! Keine unnötigen Annahmen!
- Die Gaußverteilungen sind Teil der "exponential family" oder der "log-linear models", d.h. man kann schreiben

$$p(\mathbf{x}) \propto \exp\left(\boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x})\right).$$

Zum Beispiel im 1-D Fall  $\boldsymbol{\theta} = \begin{pmatrix} \frac{1}{\sigma^2} \\ \frac{\mu}{\sigma^2} \end{pmatrix}$  und  $\boldsymbol{\phi}(\mathbf{x}) = \begin{pmatrix} x^2 \\ x \end{pmatrix}$ .

Die Tatsache, dass sich diese Verteilungen linear parameterisieren lassen, ist in vielen Anwendungen in der Statistik/im Machine Learning extrem wichtig und hilfreich.

All diese Eigenschaften machen die Gaußverteilungen zu einer flexiblen, ausdrucksstarken und leicht handhabbaren Klasse von Verteilungen. Sie werden daher sehr gerne im Machine Learning verwendet. (Häretiker könnten sagen, dass derjenige, der die Gaußverteilungen vollständig verstanden hat, auch das Machine Learning komplett beherrscht ;)

## 2 Grundzüge des Unüberwachten Lernens

Englisch: Unsupervised Learning

**Problem 2.1.** Gegeben seien Datenpunkte  $\mathbf{x}_i \in \mathbb{R}^M$ ,  $i = 1, \dots, N$  (z.B. digitalisierte Bilder). Wie kann ich diese Daten effizient und intuitiv beschreiben?

Ein paar einfache Modelle:

- Mittelwert + Noise/Rauschen:

$$\mathbf{x}_i = \boldsymbol{\mu} + \boldsymbol{\epsilon}_i \quad (7)$$

- mehrere Mittelwerte + Noise/Rauschen

$$\mathbf{x}_i = \sum_{j=1}^r \delta_{z_j, j} \boldsymbol{\mu}_j + \boldsymbol{\epsilon}_i, \quad z_j \in \{1, \dots, r\} \quad (8)$$

Dieses Modell entspricht **Clustering** und wäre genügend Stoff für eine eigene Vorlesung. Einige Stichworte sind k-Means, Bayesian k-Means, spectral clustering, hierarchical clustering, ...

- linear Modelle

$$\mathbf{x}_i = \boldsymbol{\mu} + \sum_{j=1}^r \mathbf{u}_j w_{ji} + \boldsymbol{\epsilon}_i \quad (9)$$

Dies führt zur heutigen Vorlesung, der **Hauptkomponentenanalyse (Englisch: Principal Component Analysis = PCA)**.

- nicht lineare Modelle. Ein sehr aktives Forschungsgebiet, bekannt als **manifold learning**. Stichworte sind z.B. kernelPCA, locally linear maps (LLE), ISOMAP, bi-directional GP, ...

## 3 Hauptkomponentenanalyse - PCA

**Beobachtung 3.1.** Das PCA-Problem, d.h. die Rekonstruktion von  $\mathbf{u}_j$  und  $\boldsymbol{\mu}$  gegeben nur die Datenpunkte  $\mathbf{x}_i$ , ist stark unterbestimmt. Die Anzahl der Unbekannten ist  $M + Mr + rN + MN$ , die Zahl der Gleichungen nur  $MN$ . Dies bedeutet, dass wir starke Annahmen machen müssen, um eine eindeutige Lösung zu erhalten!

Annahmen zur Herleitung der PCA:

- $\boldsymbol{\epsilon}_i$  Gaußverteilt:

$$\boldsymbol{\epsilon}_i \sim N(\boldsymbol{\epsilon}_i; \mathbf{0}, \sigma_0^2 \mathbf{1})$$

Außerdem nehmen wir an, dass  $\boldsymbol{\epsilon}_i$  klein ist. D.h.  $\sigma_0^2$  klein.

- $w_{ji}$  Gaußverteilt,

$$w_{ji} \sim N(w_{ji}; 0, \sigma_j^2)$$

Dies bedeutet dass  $w_{ji}$  unabhängig (**unkorreliert**) von  $w_{j'i}$  für  $j' \neq j$ !

Annahme zur Notation:  $\sigma_1 > \sigma_2 > \dots > \sigma_r$ .

- Sei  $\mathbf{U} = [\mathbf{u}_1 \dots \mathbf{u}_r]$ . Wir nehmen an, dass  $\mathbf{U}$  orthogonal, d.h.  $\mathbf{u}_i \perp \mathbf{u}_j$  für  $j \neq i$  und  $\|\mathbf{u}_j\| = 1$  für  $j = 1, \dots, r$ .

Herleitung der PCA:

- Neue, äquivalente Notation:

$$\mathbf{x} = \boldsymbol{\mu} + \mathbf{U}\mathbf{w} + \boldsymbol{\epsilon}$$

Hier ist  $\mathbf{U} \in \mathbb{R}^{M \times r}$ ,  $\boldsymbol{\mu} \in \mathbb{R}^M$  konstant, und  $\mathbf{w}$  zufällig verteilt wie  $\mathbf{w} \sim N(\mathbf{w}; 0, \mathbf{D})$  mit  $\mathbf{D} = \text{diag}(\sigma_1^2, \dots, \sigma_r^2)$ . Außerdem ist  $\boldsymbol{\epsilon} \in \mathbb{R}^M$  verteilt nach  $\boldsymbol{\epsilon} \sim N(\boldsymbol{\epsilon}; 0, \sigma_0^2 \mathbf{1})$ . Damit ist

$$\mathbf{x} \sim N(\mathbf{x}; \boldsymbol{\mu}, \mathbf{U}\mathbf{D}\mathbf{U}^T + \sigma_0^2 \mathbf{1}).$$

- Nun ergänze  $\mathbf{U}$  zu einer orthonormalen Basis  $\bar{\mathbf{U}} \in \mathbb{R}^{M \times M}$ , d.h.  $\bar{\mathbf{U}}^T \bar{\mathbf{U}} = \bar{\mathbf{U}} \bar{\mathbf{U}}^T = \mathbf{1}$  und

$$\begin{aligned} \mathbf{x} &\sim N(\mathbf{x}; \boldsymbol{\mu}, \mathbf{U}\mathbf{D}\mathbf{U}^T + \sigma_0^2 \mathbf{1}) \\ &\sim N(\mathbf{x}; \boldsymbol{\mu}, \mathbf{U}\mathbf{D}\mathbf{U}^T + \sigma_0^2 \bar{\mathbf{U}} \bar{\mathbf{U}}^T) \\ &\sim N(\mathbf{x}; \boldsymbol{\mu}, \underbrace{\bar{\mathbf{U}} \bar{\mathbf{D}} \bar{\mathbf{U}}^T}_{=\boldsymbol{\Sigma}}), \end{aligned}$$

wobei  $\bar{\mathbf{D}} = \text{diag}(\sigma_1^2 + \sigma_0^2, \dots, \sigma_r^2 + \sigma_0^2, \sigma_0^2, \dots, \sigma_0^2)$ .

**Idee PCA:** Schätze empirische Kovarianz  $\hat{\boldsymbol{\Sigma}}$  aus den gegebenen Daten und berechne Eigenwertzerlegung. Die Eigenvektoren entsprechen den gesuchten Basisvektoren!

**Kochrezept PCA:**

1. Schätze Mittelwert und normalisiere Daten auf Mittelwert 0.

$$\hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_i \mathbf{x}_i, \quad \tilde{\mathbf{x}}_i = \mathbf{x}_i - \hat{\boldsymbol{\mu}}$$

(Manchmal, z.B. bei LSA, wird dieser Schritt weggelassen.)

2. Schätze Kovarianzmatrix  $\hat{\boldsymbol{\Sigma}}$

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{N} \sum_i \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T$$

3. Berechne Eigenwertzerlegung von  $\hat{\boldsymbol{\Sigma}}$

$$\hat{\boldsymbol{\Sigma}} = \hat{\mathbf{U}} \hat{\mathbf{D}} \hat{\mathbf{U}}^T$$

wobei  $\hat{\mathbf{D}} = \text{diag}(\hat{\lambda}_1, \dots, \hat{\lambda}_M)$ . Permutiere die Eigenwertzerlegung so, dass  $\hat{\lambda}_1 > \dots > \hat{\lambda}_M$ .

4. Bestimme  $\hat{r}$  aus dem Abfall des Eigenwertspektrums (Zur Erinnerung:  $\sigma_0^2 < \sigma_1^2, \dots, \sigma_{\hat{r}}^2$ )

Die ersten  $\hat{r}$  Spalten von  $\hat{\mathbf{U}}$  sind dann die gesuchten Basisvektoren  $\hat{\mathbf{U}} = [\hat{\mathbf{u}}_1, \dots, \hat{\mathbf{u}}_{\hat{r}}]$ . Die Eigenwerte  $\hat{\lambda}_j$  sind Schätzungen der Varianz in der jeweiligen Richtung  $\hat{\sigma}_j^2 = \hat{\lambda}_j$ . Das Abschneiden des Eigenwertspektrums durch eine geeignete Wahl von  $\hat{r}$  begründet dann auch den Namen "Hauptkomponentenanalyse".

Beweisskizze PCA: Die Eigenwertzerlegung ist eindeutig, falls alle Eigenwerte verschieden und richtig sortiert.

Bemerkungen zur PCA:

- Für einen neuen Datenpunkt  $\mathbf{z} \in \mathbb{R}^M$  erhalten wir die Hauptkomponenten  $\mathbf{w}$  einfach durch

$$\mathbf{w} = \hat{\mathbf{U}}^T (\mathbf{z} - \hat{\boldsymbol{\mu}}).$$

- Effiziente Berechnung der PCA: Die oben vorgestellte Methode skaliert  $O(M^3)$  wegen der Eigenwertzerlegung. Für große  $M$  ist dies unpraktisch. Man kann die Eigenwertzerlegung aber mit Hilfe einer **Singulärwertzerlegung** umgehen. Falls  $N \ll M$  ist dies viel effizienter, skaliert nämlich mit  $O(MN^2)$ .

Sei  $\mathbf{X} = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_N]$ . Die Schätzung der Kovarianz kann dann auch geschrieben werden als  $\hat{\boldsymbol{\Sigma}} = \frac{1}{N} \mathbf{X} \mathbf{X}^T$ . Jetzt besitzt aber jede Matrix eine Singulärwertzerlegung. D.h. für  $\mathbf{X}$  können wir schreiben

$$\mathbf{X} = \mathbf{U}' \mathbf{S} \mathbf{V}',$$

wobei  $\mathbf{U}' \in \mathbb{R}^{N \times N}$  und  $\mathbf{V}' \in \mathbb{R}^{M \times M}$  unitär/orthogonal und  $\mathbf{S} \in \mathbb{R}^{M \times N}$  diagonal ist, d.h.  $\mathbf{S} = \text{diag}(\sigma_1, \dots, \sigma_{\min(N,M)})$ . Bisher haben wir berechnet

$$\hat{\Sigma} = \frac{1}{N} \mathbf{X} \mathbf{X}^T = \hat{\mathbf{U}} \hat{\mathbf{D}} \hat{\mathbf{U}}^T,$$

mit der Singulärwertzerlegung erhalten wir

$$\frac{1}{N} \mathbf{X} \mathbf{X}^T = \frac{1}{N} \mathbf{U}' \mathbf{S} \mathbf{V}'^T \mathbf{S}^T \mathbf{V}'^T = \frac{1}{N} \mathbf{U}' \underbrace{\mathbf{S} \mathbf{S}^T}_{=\mathbf{D}'} \mathbf{V}'^T,$$

wobei  $\mathbf{D}' = \text{diag}(\sigma_1^2, \dots, \sigma_{\min(N,M)}^2, 0, \dots, 0)$  diagonal ist.

Zusammen bedeutet dies, dass wir bei der PCA aus einer kostengünstig berechenbaren Singulärwertzerlegung dieselben Ergebnisse wie aus einer Eigenwertzerlegung erhalten können.

(Mann kann sogar  $M \rightarrow \infty$  gehen lassen, siehe kernelPCA!)

- Wir haben die PCA unter der Annahme eines Gaußverteilten linearen Modells hergeleitet. Oft wird PCA aber auch dann verwandt, wenn die **Modellannahmen verletzt** sind oder gar **kein Modell bekannt** ist. Oft liefert dies trotzdem gute Ergebnisse, siehe der experimentelle Teil der Vorlesung.

Der wichtige Grundgedanke bei der PCA ist, dass die Datenpunkte auf einer Linie/Fläche (einem linearen Unterraum des  $\mathbb{R}^M$ ) liegen und dass wir diesen aus der Kovarianz der Daten schätzen können.