



Grundlagen der Künstlichen Intelligenz

Wintersemester 25/26

Vorlesung 8

Expectation-Maximization, Gaussian Mixture Models und Variational Autoencoder

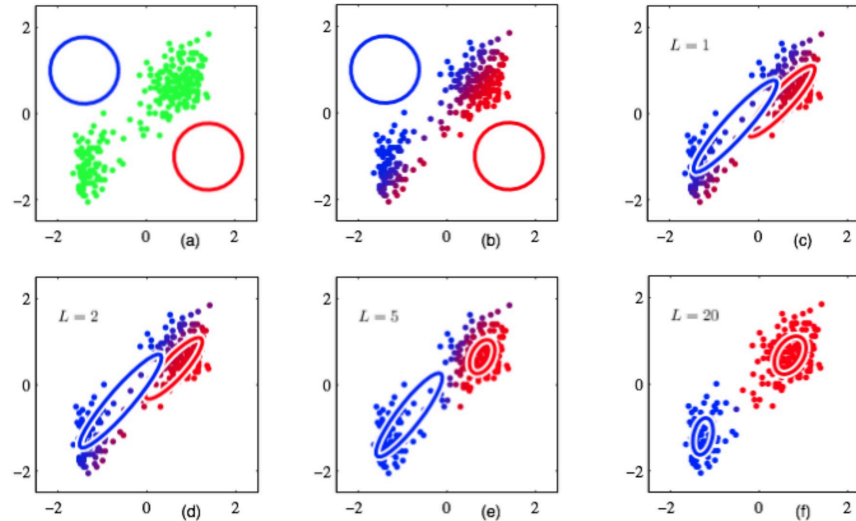
Prof. Dr. Pascal Friederich
T.T.-Prof. Dr. Peer Nowack

Lernergebnisse

Was werden wir heute lernen?

- Gaussian Mixture Models und Modelle mit latenten Variablen
- Expectation-Maximization Algorithmus (EM) für das Training von Mixture Models
- Herleitung von EM als Lösungsalgorithmus für allgemeine Modelle mit latenten Variablen
- Analyse des EM-Algorithmus
- Kurzer Einblick in Variational Autoencoder und deren Beziehung zu EM

Teil 1



1) Expectation Maximization und Gaussian Mixture Models

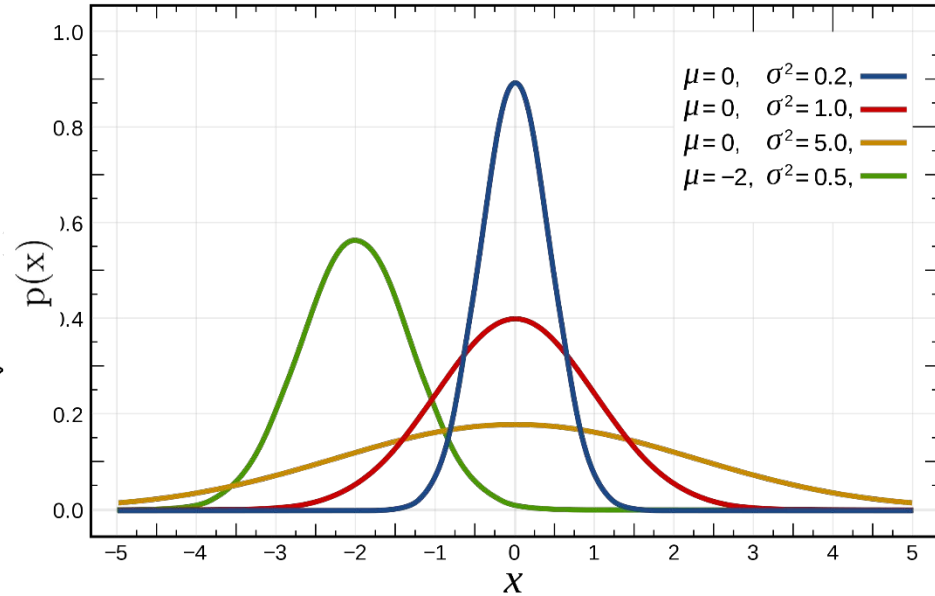
2) Verallgemeinerte Sicht auf EM

3) Variational Autoencoder

Gauss Verteilung

- Kontinuierliche Zufallsvariable: $X \in \mathbb{R}$
- Distribution is completely specified by mean μ and variance σ^2

$$p(x) = \mathcal{N}(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x - \mu)^2}{2\sigma^2}\right\}$$



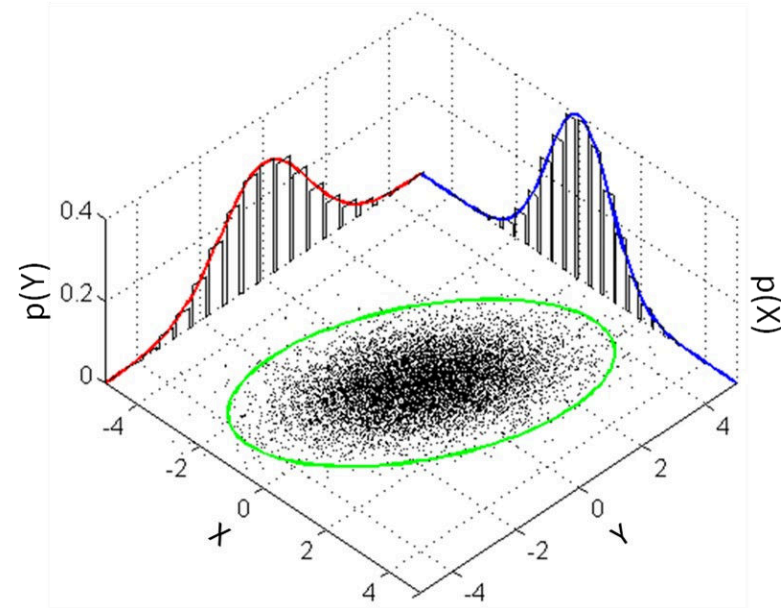
Multivariate Gauss Verteilung

- Kontinuierliche Zufallsvariable :

$$X \in \mathbb{R}^d$$

- Verteilung ist spezifiziert durch Mittelwertvektor $\boldsymbol{\mu}$ und Kovarianz-Matrix $\boldsymbol{\Sigma}$

$$p(\boldsymbol{x}) = \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{|2\pi\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{(\boldsymbol{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{x} - \boldsymbol{\mu})}{2} \right\}$$



Mixture Modelle

Eine Mixture-Verteilung ist die **Summe der Einzelverteilungen**:

$$p(\mathbf{x}) = \sum_{k=1}^K p(k) p(\mathbf{x}|k)$$

Anzahl der Komponenten

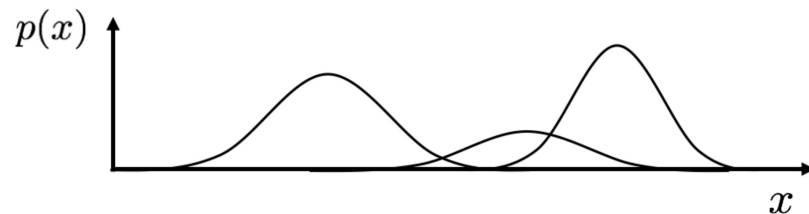
Mixture Koeffizienten

k-te Mixture Komponente

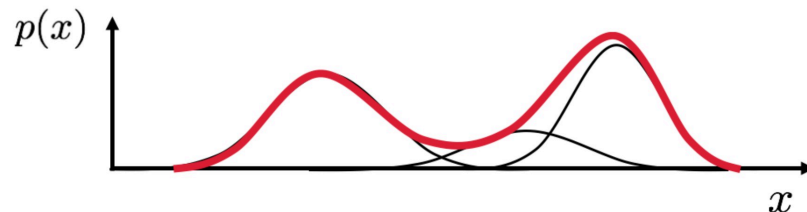
Im Grenzfall mit vielen/unendlichen Komponenten kann dies **eine beliebige glatte Dichte approximieren**

Gaussian Mixture Models (GMM)

- Einzelne Gauß Verteilungen



- Summe der Gaußschen Verteilungen



Gaußsche Mixture Modelle (GMMs)

- **Mixture Koeffizienten:**

$$p(k) = \pi_k, \text{ with } 0 \leq \pi_k \leq 1, \sum_k \pi_k = 1$$

- **Mixture Komponente:**

$$p(\mathbf{x}|k) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- **Mixture Verteilung:**

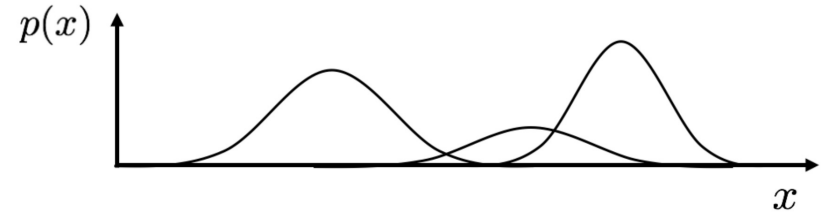
$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- Integriert immer zu 1
- Parameter der Mixture

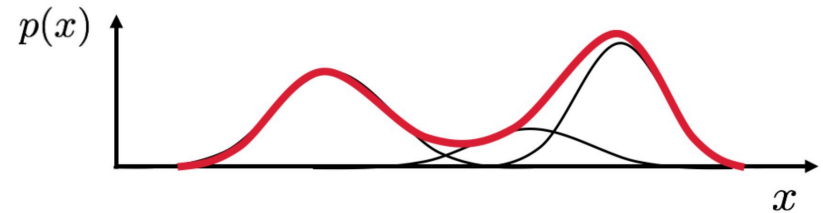
$$\boldsymbol{\theta} = \{\pi_1, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \dots, \pi_K, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K\}$$

Gaussian Mixture Models (GMM)

- Einzelne Gauß Verteilungen



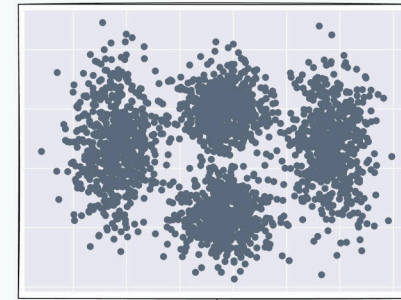
- Summe der Gaußschen Verteilungen



GMMs als Clustering Methode

Gaussian Mixture Model vs. KMeans

 blog.DailyDoseofDS.com



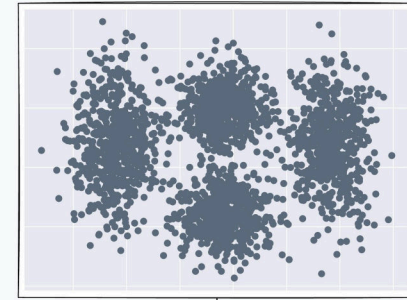
Clusters with
different
variances

GMMs als Clustering Methode

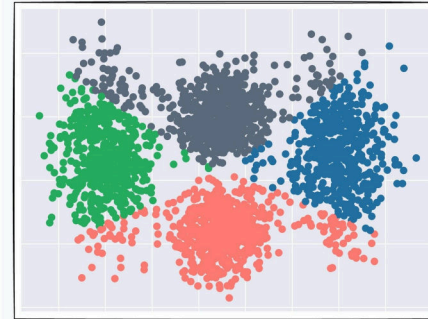
- Die Kovarianz Matrix lernt die (elliptische) „Form“ des Clusters
- K-Means: „Form“ ist fixiert durch kürzesten Abstand

Gaussian Mixture Model vs. KMeans

 blog.DailyDoseofDS.com



KMeans



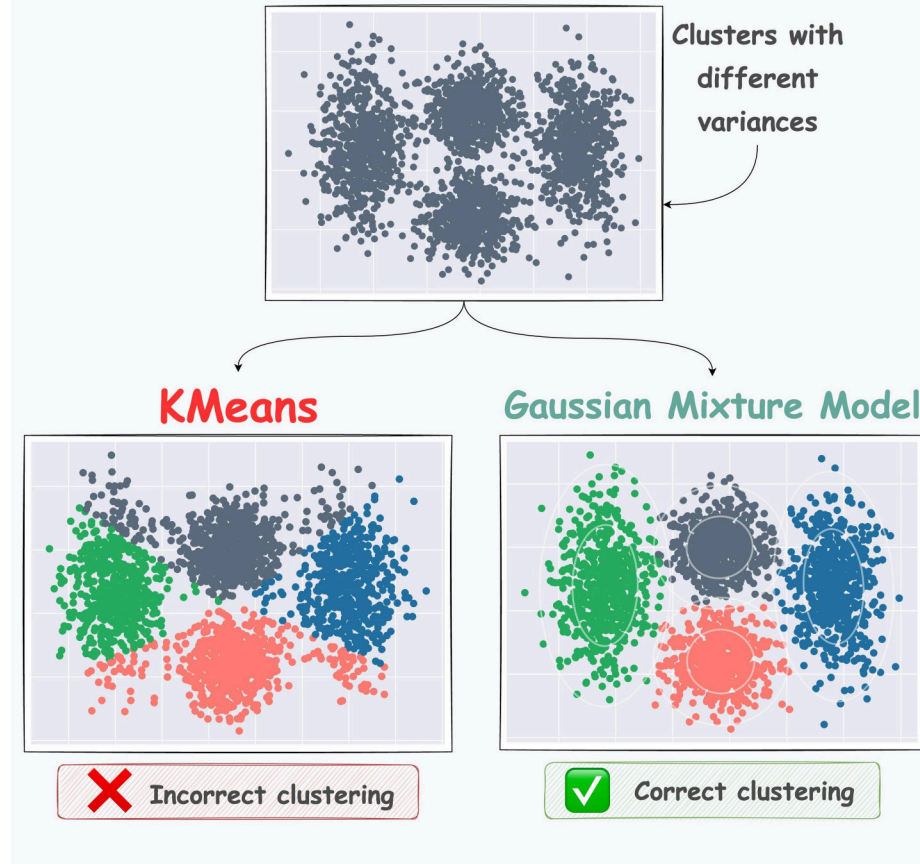
Incorrect clustering

GMMs als Clustering Methode

- Die Kovarianz Matrix lernt die (elliptische) „Form“ des Clusters
- K-Means: „Form“ ist fixiert durch kürzesten Abstand

Gaussian Mixture Model vs. KMeans

 blog.DailyDoseofDS.com



Maximum Likelihood einer Mixture Verteilung

- (Marginal-)Log-Likelihood mit N iid. Punkten

$$\mathcal{L} = \log L(\boldsymbol{\theta}) = \sum_{i=1}^N \log \underbrace{p(\mathbf{x}_i | \boldsymbol{\theta})}_{\text{marginal}} = \sum_{i=1}^N \log \underbrace{\left(\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)}_{\text{non-exponential family}}$$

- F: Können wir einen Gradientenaufstieg durchführen?

(iid: Independent and identically distributed)

Maximum Likelihood einer Mixture Verteilung

- (Marginal-)Log-Likelihood mit N iid. Punkten

$$\mathcal{L} = \log L(\boldsymbol{\theta}) = \sum_{i=1}^N \log \underbrace{p(\mathbf{x}_i | \boldsymbol{\theta})}_{\text{marginal}} = \sum_{i=1}^N \log \underbrace{\left(\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)}_{\text{non-exponential family}}$$

- F: Können wir einen Gradientenaufstieg durchführen?

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \boldsymbol{\mu}_j} &= \sum_{i=1}^N \frac{\pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j) \\ &= \sum_{i=1}^N \frac{p(j)p(\mathbf{x}_i | j)}{p(\mathbf{x}_i)} \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j) \\ &= \sum_{i=1}^N \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j) p(j | \mathbf{x}_i) \end{aligned}$$

(iid: Independent and identically distributed)

Maximum Likelihood einer Mixture Verteilung

- (Marginal-)Log-Likelihood mit N iid. Punkten

$$\mathcal{L} = \log L(\boldsymbol{\theta}) = \sum_{i=1}^N \log \underbrace{p(\mathbf{x}_i | \boldsymbol{\theta})}_{\text{marginal}} = \sum_{i=1}^N \log \underbrace{\left(\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)}_{\text{non-exponential family}}$$

- **F:** Können wir einen **Gradientenaufstieg** durchführen?
 - × Der Gradient hängt von allen anderen Komponenten ab (zyklische Abhängigkeit)
 - × Keine geschlossene Lösung
 - × Typischerweise sehr langsame Konvergenz
 - ▶ **A:** Ja, aber die **Summe** (Marginalisierung) **verträgt sich nicht gut mit dem Log**

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \boldsymbol{\mu}_j} &= \sum_{i=1}^N \frac{\pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j) \\ &= \sum_{i=1}^N \frac{p(j) p(\mathbf{x}_i | j)}{p(\mathbf{x}_i)} \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j) \\ &= \sum_{i=1}^N \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j) p(j | \mathbf{x}_i) \end{aligned}$$

(iid: Independent and identically distributed)

Maximum Likelihood einer Mixture Verteilung

Gradient Descent ist möglich... aber sehr langsam!

- Können wir ein spezielleres Optimierungsverfahren für Mixture Modelle finden?
- Ja! **Expectation Maximization!**

Schätzung von Gaußschen Mixture Modellen

Warum also ist die Optimierung von $\mathcal{L} = \sum_{i=1}^N \log p_{\theta}(\mathbf{x}_i)$ so schwierig?

- Wir wissen nicht, welche Mixture-Komponente k welchen Datenpunkt erzeugt hat
- Wenn wir **Daten aus der gemeinsamen Verteilung** $p(\mathbf{x}_i, k_i | \theta)$ hätten, dann wäre es einfach ...

In diesem Fall können wir einfach eine Maximum-Likelihood-Schätzung durchführen:

- Zuweisungen: $q_{ik} = p(k_i = k | \mathbf{x}_i)$ binär: $q_{ik} = \mathbb{I}(k, k_i)$ hat Wert 1, wenn die i -te Probe zur k -ten Komponente gehört, und 0, wenn nicht
- Koeffizienten: $\pi_k = \frac{\sum_i q_{ik}}{N}$
- Mittelwerte: $\boldsymbol{\mu}_k = \frac{\sum_i q_{ik} \mathbf{x}_i}{\sum_i q_{ik}}$
- Kovarianzen: $\boldsymbol{\Sigma}_k = \frac{\sum_i q_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T}{\sum_i q_{ik}}$

Expectation-Maximization: Henne Ei...

Wir wissen jedoch nicht, welche Komponente zu welcher Probe gehört.

- Können wir das abschätzen? **Bei einem gegebenen Mixture-Modell, ja!**

Expectation-Maximization: Henne Ei...

Wir wissen jedoch nicht, welche Komponente zu welcher Probe gehört.

- Können wir das abschätzen? **Bei einem gegebenen Mixture-Modell, ja!**

Expectation-Step:

- Berechnung der Clusterwahrscheinlichkeiten aka **Responsibilities** für jedes Sample (Bayes-Regel)

$$q_{ik} = p(k_i = k | \mathbf{x}_i) = \frac{p(\mathbf{x}_i | k)p(k)}{\sum_j p(\mathbf{x}_i | j)p(j)} = \frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

- Die Responsibilities q_{ik} liegen zwischen 0 und 1
- Aber wir müssen **die Gaußschen Komponenten kennen**

Expectation-Maximization: Henne Ei...

Wir wissen jedoch nicht, welche Komponente zu welcher Probe gehört.

- Können wir das abschätzen? **Bei einem aktuellen Mixture-Modell, ja!**

Maximization-Step:

- Berechnung der (gewichteten) Maximum-Likelihood-Schätzung

$$\pi_k = \frac{\sum_i q_{ik}}{N}$$

$$\boldsymbol{\mu}_k = \frac{\sum_i q_{ik} \boldsymbol{x}_i}{\sum_i q_{ik}}$$

$$\boldsymbol{\Sigma}_k = \frac{\sum_i q_{ik} (\boldsymbol{x}_i - \boldsymbol{\mu}_k)(\boldsymbol{x}_i - \boldsymbol{\mu}_k)^T}{\sum_i q_{ik}}$$

- Aber wir müssen **die Responsibilities kennen**

Algorithmus: EM für GMMs

- **Initialisieren:** Mixture-Komponenten + Mixture-Koeffizienten
 - z.B. Verwendung von k-means für die Komponentenmittelwerte und eine anfängliche Kovarianz
- **Wiederhole den Vorgang** bis zur Konvergenz:
 - **Expectation-Step:** Berechnung der Responsibilities

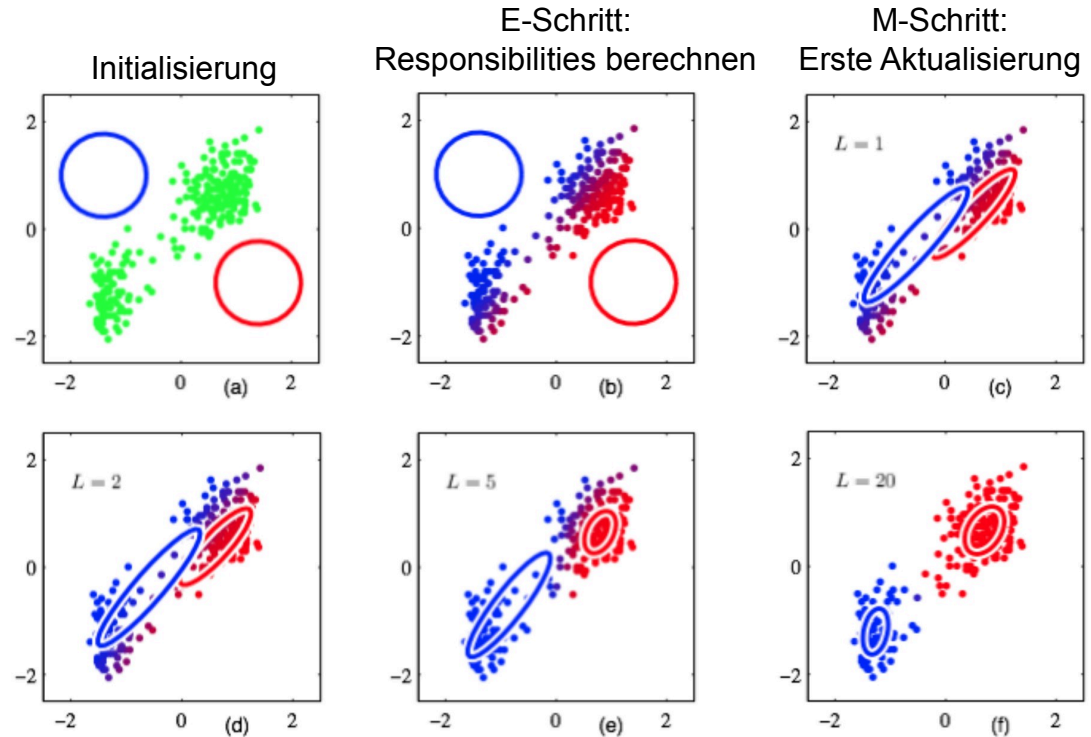
$$q_{ik} = \frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

- **Maximization-Step:** Aktualisierung der Koeffizienten, Komponentenmittelwerte und Komponentenvarianz

$$\pi_k = \frac{\sum_i q_{ik}}{N} \quad \boldsymbol{\mu}_k = \frac{\sum_i q_{ik} \mathbf{x}_i}{\sum_i q_{ik}} \quad \boldsymbol{\Sigma}_k = \frac{\sum_i q_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T}{\sum_i q_{ik}}$$

Abbildung

- Jede Komponente steht für ein Cluster im Datensatz
- EM ist sehr empfindlich gegenüber der Initialisierung



Praktische Überlegungen...

Wie **viele Komponenten** brauchen wir?

- Mehr Komponenten führen in der Regel zu einem besseren Trainings-Likelihood
- Aber sind mehr Komponenten unbedingt besser? Nein, wegen der **Overfitting!**
- Es handelt sich wieder um ein **Model-Selection Problem** (verwende z.B. Cross-Validation)

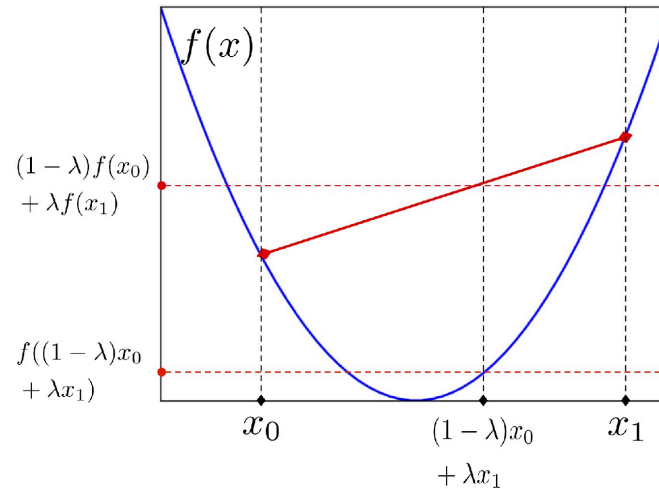
EM versus k-means

K-means kann als Spezialfall von EM angesehen werden:

- Co-Varianzen werden immer auf 0 gesetzt (im Grenzfall)
- **E-Step / Zuweisungsschritt:**
 - Responsibilities q_{ik} des nächstgelegenen Clusters k werden auf 1 gesetzt, alle anderen Werte sind 0
- **M-Step / Anpassungsschritt:**
 - Die Aktualisierung für den Mittelwert ist die gleiche
 - Co-Varianzen werden ignoriert (auf nahe 0 gesetzt)
- EM ist schwieriger zu erlernen als k-means, aber es liefert auch Varianzen und Dichten
- Häufig wird k-means verwendet, um die Mittelwerte für EM zu initialisieren.

K-means ist dafür bekannt, dass es konvergiert. Konvergiert auch EM immer?

Teil 2



1) Expectation Maximization und Gaussian Mixture Models

2) Verallgemeinerte Sicht auf EM

3) Variational Autoencoder

Modelle für latente Variablen

Mixture-Modelle sind ein Beispiel für Modelle **mit latenten Variablen**

- Beobachtete Variablen: \mathbf{x} , latente Variablen: \mathbf{z} (z. B. der Index der Mixture-Komponente)

- **Parametrisches Modell:** $p_{\theta}(\mathbf{x}, \mathbf{z})$

- **Marginale Verteilung:** $p_{\theta}(\mathbf{x}) = \underbrace{\sum_z p(\mathbf{x}, \mathbf{z})}_{\text{discrete latent variable}}, \quad p_{\theta}(\mathbf{x}) = \underbrace{\int_{\mathbf{z}} p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z}}_{\text{continuous latent variable}}$

(Marginal) Log-Likelihood:

$$L(\theta) = \sum_{i=1}^N \log p_{\theta}(\mathbf{x}_i) = \sum_{i=1}^N \log \left(\sum_{z_i} p_{\theta}(\mathbf{x}_i, z_i) \right)$$

... ist für alle Modelle mit latenten Variablen **schwer zu optimieren** (wegen des Logarithmus einer Summe)

Expectation-Maximization (EM)

Expectation-Maximization (EM) ist ein **allgemeiner Algorithmus** zur Schätzung von **Modellen mit latenten Variablen**

- **Häufigste Anwendung:** Gaußsche Mixture Modelle
- ... aber auch viele andere (tiefe) Modelle
- Seine Erweiterung wird **Variational Bayes** genannt, die dem Variational Autoencoder zugrunde liegt

EM kann auf 2 Arten abgeleitet werden:

- Jensen's Ungleichung
- Zerlegung in Untergrenze und Kullback-Leibler-Divergenz-Begriff
(besserer, aber komplizierterer Weg... siehe ML-Vorlesung)

Jensens Ungleichung: Erinnerung

- Wenn eine Funktion konvex ist, dann

$$f\left(\sum_i \lambda_i \mathbf{x}_i\right) \leq \sum_i \lambda_i f(\mathbf{x}_i),$$

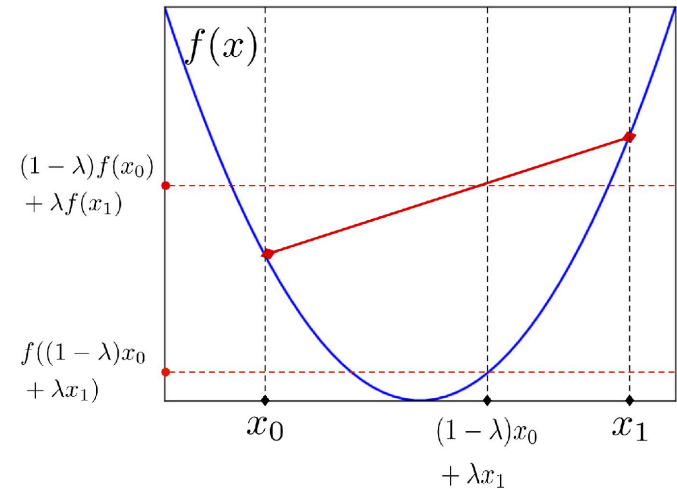
wobei $\{\lambda_i\}$ so sind, dass jedes $\lambda_i \geq 0$ und

$$\sum_i \lambda_i = 1$$

- Wenn wir die λ_i als die Parameter einer kategorialen Verteilung behandeln, d. h. $\lambda_i = p(\mathbf{x} = \mathbf{x}_i)$, kann dies wie folgt geschrieben werden

$$f(\mathbb{E}_p[\mathbf{x}]) \leq \mathbb{E}_p[f(\mathbf{x})].$$

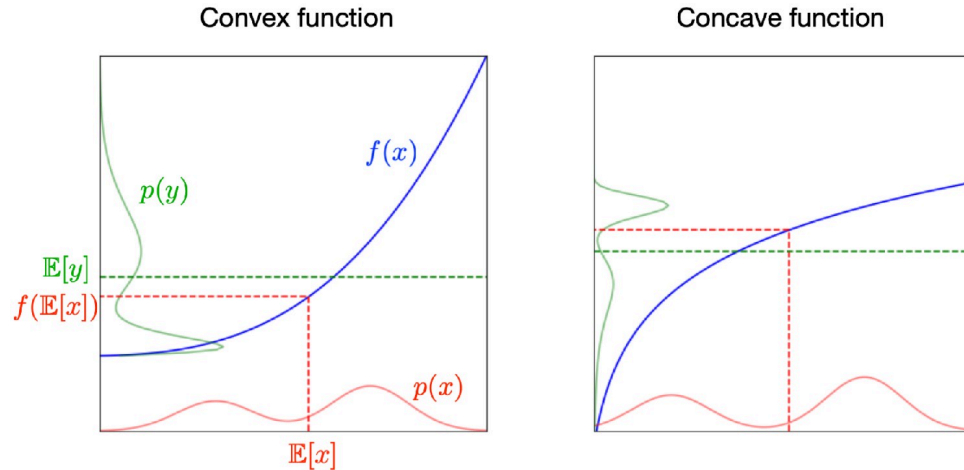
- Dies ist als **Jensensche Ungleichung** bekannt. Sie gilt auch für kontinuierliche Verteilungen.



Jensensche Inequality

- Eine Funktion $f(\mathbf{x})$ ist **konkav**, wenn $-f(\mathbf{x})$ konvex ist. In diesem Fall drehen wir die Jensensche Inequality um:

$$f(\mathbb{E}_p[\mathbf{x}]) \geq \mathbb{E}_p[f(\mathbf{x})].$$



Zurück zu Expectation-Maximization

(Marginale) Log-Likelihood:

$$L(\boldsymbol{\theta}) = \sum_{i=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{x}_i) = \sum_{i=1}^N \log \left(\sum_{z_i} p_{\boldsymbol{\theta}}(\mathbf{x}_i, z_i) \right)$$

Konkave Funktion!

... die für alle Modelle mit latenten Variablen **schwer zu optimieren** ist (wegen des Logarithmus einer Summe)

Zurück zu Expectation-Maximization

(Marginale) Log-Likelihood:

$$L(\boldsymbol{\theta}) = \sum_{i=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{x}_i) = \sum_{i=1}^N \log \left(\sum_{z_i} p_{\boldsymbol{\theta}}(\mathbf{x}_i, z_i) \right)$$

Konkave Funktion!

... die für alle Modelle mit latenten Variablen **schwer zu optimieren** ist (wegen des Logarithmus einer Summe)

- Trick: Führe neue Verteilung $q_i(z_i)$ ein (wir werden gleich sehen, was das ist):

$$L(\boldsymbol{\theta}) = \sum_{i=1}^N \log \left(\sum_{z_i} q_i(z_i) \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_i, z_i)}{q_i(z_i)} \right) = \sum_{i=1}^N \log \mathbb{E}_{q_i(z_i)} \left[\frac{p_{\boldsymbol{\theta}}(\mathbf{x}_i, z_i)}{q_i(z_i)} \right]$$

Zurück zu Expectation-Maximization

(Marginale) Log-Likelihood:

$$L(\boldsymbol{\theta}) = \sum_{i=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{x}_i) = \sum_{i=1}^N \log \left(\sum_{z_i} p_{\boldsymbol{\theta}}(\mathbf{x}_i, z_i) \right)$$

Konkave Funktion!

... die für alle Modelle mit latenten Variablen **schwer zu optimieren** ist (wegen des Logarithmus einer Summe)

- Trick: Führe neue Verteilung $q_i(z_i)$ ein (wir werden gleich sehen, was das ist):

$$L(\boldsymbol{\theta}) = \sum_{i=1}^N \log \left(\sum_{z_i} q_i(z_i) \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_i, z_i)}{q_i(z_i)} \right) = \sum_{i=1}^N \log \mathbb{E}_{q_i(z_i)} \left[\frac{p_{\boldsymbol{\theta}}(\mathbf{x}_i, z_i)}{q_i(z_i)} \right]$$

- Log ist eine konkave Funktion ist. Daher können wir die Jensensche Ungleichung verwenden, um **den Logarithmus in den Erwartungswert zu schieben** und so die **untere Schranke** zu erhalten:

$$\log(\mathbb{E}_p[\mathbf{x}]) \geq \mathbb{E}_p[\log(\mathbf{x})]$$

$$L(\boldsymbol{\theta}) \geq \sum_{i=1}^N \mathbb{E}_{q_i(z_i)} \left[\log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_i, z_i)}{q_i(z_i)} \right] := \mathcal{L}(q, \boldsymbol{\theta})$$

Untere Variationsschranke

- Soeben wurde eine untere Schranke für die LogLikelihood abgeleitet

$$L(\boldsymbol{\theta}) \geq \sum_{i=1}^N \mathbb{E}_{q_i(z_i)} \left[\log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_i, z_i)}{q_i(z_i)} \right] := \mathcal{L}(q, \boldsymbol{\theta})$$

- Vereinfachen:

$$\mathcal{L}(q, \boldsymbol{\theta}) = \sum_{i=1}^N \mathbb{E}_{q_i(z_i)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_i, z_i)] - \underbrace{\mathbb{E}_{q_i(z_i)} [\log q_i(z_i)]}_{\text{constant w.r.t } \boldsymbol{\theta}}$$

Anmerkungen:

- Die untere Schranke ist viel einfacher zu optimieren
- Oft gibt es geschlossene Formeln für die Herleitung von q gegeben $\boldsymbol{\theta}$ und umgekehrt
- Die untere Grenze **gilt für jede Wahl** von $q_i(z_i)$

Expectation-Maximization Schritte

EM optimiert die untere Schranke durch die Anwendung von 2 Schritten:

- **(M)aximization-step** (finde Maximum von blau):

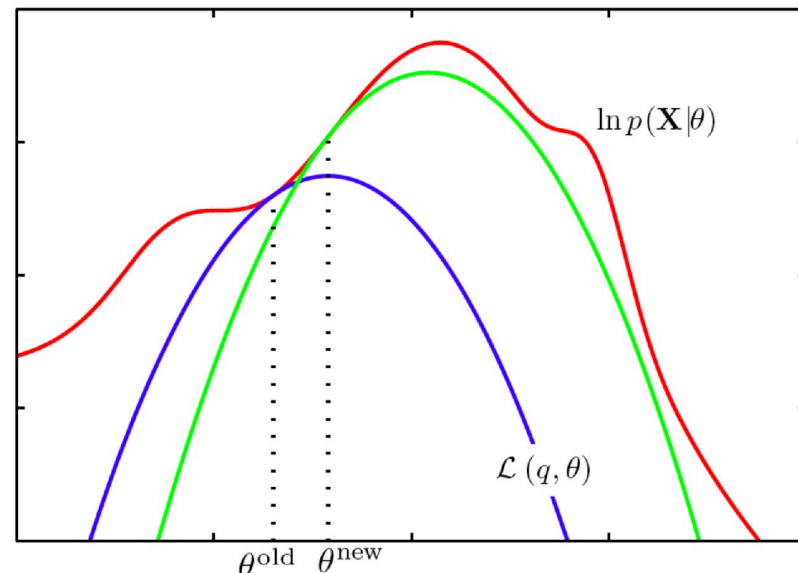
- Fixe Variationsverteilung $q = q_{\text{old}}$
- Maximierung der unteren Schranke in Bezug auf θ

$$\theta_{\text{new}} = \operatorname{argmax}_{\theta} \mathcal{L}(q_{\text{old}}, \theta)$$

- **(E)xpectation-step** (blau \rightarrow grün):

- Parameter festlegen $\theta = \theta_{\text{new}}$
- Maximierung der unteren Schranke in Bezug auf q

$$q_{\text{new}} = \operatorname{argmax}_q \mathcal{L}(q, \theta_{\text{new}})$$



Expectation-Maximization Schritte

Expectation-Step:

- Finde $q(z)$ welches die untere Schranke maximiert:

$$q_{\text{new}} = \operatorname{argmax}_q \mathcal{L}(q, \boldsymbol{\theta}_{\text{old}})$$

- Kann in **geschlossener Form für diskrete latent Variablen** z durchgeführt werden:
 - Z.B. Komponentenindizes für Mixture-Modelle
- Das optimale q entspricht dem Posterior (hier nicht hergeleitet)

$$q_{\text{new}}(z_i) = p_{\boldsymbol{\theta}_{\text{old}}}(z_i | \mathbf{x}_i) = \frac{p_{\boldsymbol{\theta}_{\text{old}}}(\mathbf{x}_i, z_i)}{p_{\boldsymbol{\theta}_{\text{old}}}(\mathbf{x}_i)}$$

Expectation-Maximization Schritte

(M)aximierungs-Schritt:

- Maximierung der unteren Schranke bezüglich θ

$$\begin{aligned}\theta &= \arg \max_{\theta} \mathcal{L}(q, \theta) = \sum_{i=1}^N \mathbb{E}_{q_i(z_i)} [\log p_{\theta}(\mathbf{x}_i, z_i)] - \underbrace{\mathbb{E}_{q_i(z_i)} [\log q_i(z_i)]}_{\text{constant w.r.t } \theta} \\ &= \arg \max_{\theta} \sum_i \sum_{z_i} q(z_i) \log p(\mathbf{x}_i, z_i | \theta) + \text{const}\end{aligned}$$

- Jeder mögliche Wert der “fehlenden Daten” wird mit der Wahrscheinlichkeit gewichtet, dass dieser Wert wahr ist, d.h.:

$$q(z_i) = p_{\theta_{\text{old}}}(z_i | \mathbf{x}_i) = \frac{p_{\theta_{\text{old}}}(\mathbf{x}_i, z_i)}{p_{\theta_{\text{old}}}(\mathbf{x}_i)}$$

... erhalten im E-Schritt

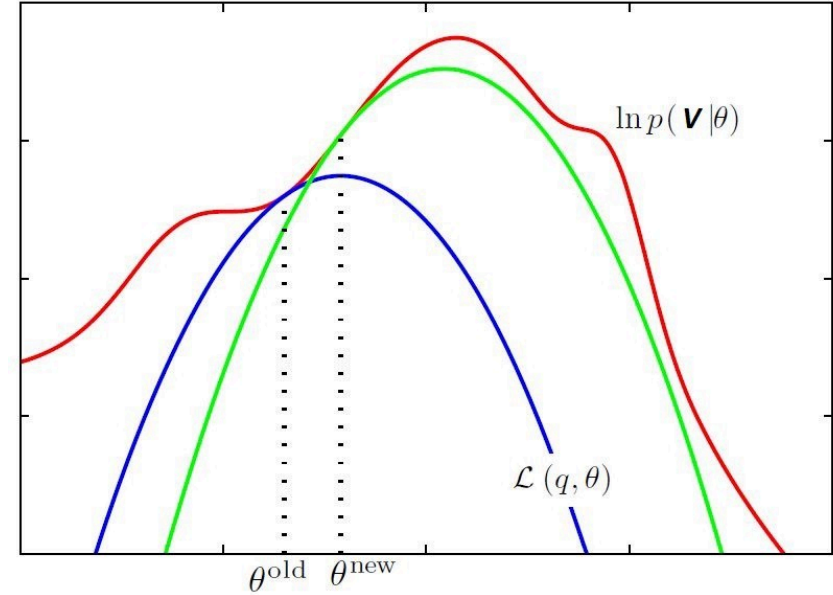
Illustration von EM

Optimierung

- Die untere Schranke (blaue Kurve) ist eine **konvexe Annäherung an** den Marginal Likelihood (rote Kurve)
 - Das Maximum der unteren Schranke kann leicht ermittelt werden (θ_{new})
 - Lösungen θ_{new} in geschlossener Form verfügbar, kein Gradientenabstieg erforderlich
- Berechne eine neue untere Schranke (grüne Kurve)
- Aufgrund der lokalen Approximation der unteren Schranke **kann EM nur lokale Optima finden**

Konvergenz: **EM konvergiert immer!**

- EM verbessert die untere Grenze bei jedem Schritt
- Es kann gezeigt werden (siehe ML-Vorlesung), dass auch die **logarithmische Grenzwahrscheinlichkeit immer erhöht wird** (oder gleich bleibt)



EM für GMMs revisited

- **Mixture Koeffizienten:**

$$p(k) = \pi_k, \text{ with } 0 \leq \pi_k \leq 1, \sum_k \pi_k = 1$$

- **Mixture Komponenten:**

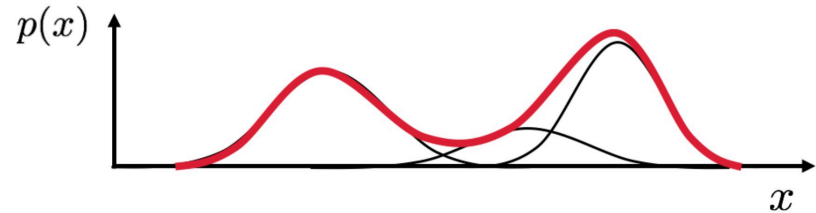
$$p(\mathbf{x}|k) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- **Mixture Verteilung:**

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- Integriert immer zu 1
- Parameter des Mixtures

$$\boldsymbol{\theta} = \{\pi_1, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \dots, \pi_K, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K\}$$



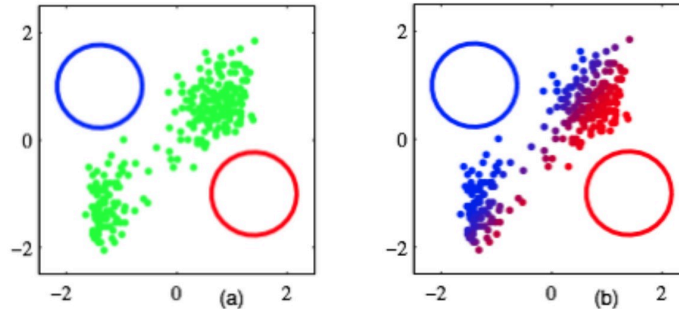
EM für GMMs revisited

E-Step

- Berechne "Responsibilities":

$$\begin{aligned} \text{Folie 33} \quad \downarrow \quad \text{Bayes} \\ q_{ik} &= p(z_i = k | \mathbf{x}_i) = \frac{p(\mathbf{x}_i | z_i = k)p(z_i = k)}{p(\mathbf{x}_i)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \end{aligned}$$

- Wie viel die Komponente k zur Erzeugung von \mathbf{x}_i beiträgt gemäß dem aktuellen Mixture-Modell



EM für GMMs revisited

M-Step:

Folie 34



$$\theta = \arg \max_{\theta} \sum_i \sum_k q_{ik} \log p(k) p(\mathbf{x}_i | k)$$

log Regeln



$$= \arg \max_{\theta} \sum_i \sum_k q_{ik} \log p(k) + \sum_k \sum_i q_{ik} \log p(\mathbf{x}_i | k)$$

- Wir können **die Aktualisierungen** der einzelnen Komponenten und Koeffizienten **trennen**
 - nur additive Terme im Lower Bound

- **Koeffizienten aktualisieren:**

$$p(k) = \pi_k$$



$$\boldsymbol{\pi} = \arg \max_{\boldsymbol{\pi}} \sum_i \sum_k q_{ik} \log \pi_k$$

- **Komponenten aktualisieren:**

$$p(\mathbf{x}|k) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$



$$\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k = \arg \max_{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k} \sum_i q_{ik} \log \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- Jeder Datenpunkt wird mit q_{ik} gewichtet
- **Gewichtete Maximum-Likelihood-Schätzung**

EM für Gaußsche Mixture-Modelle

Gewichtete Maximum-Likelihood Updates:

- Aktualisierungskoeffizienten: $\pi = \arg \max_{\pi} \sum_i \sum_k q_{ik} \log \pi_k$

EM für Gaußsche Mixture-Modelle

Gewichtete Maximum-Likelihood Updates:

- Koeffizienten aktualisieren: $\pi = \arg \max_{\pi} \sum_i \sum_k q_{ik} \log \pi_k$
 - Ergebnis: $\pi_k = \frac{\sum_i q_{ik}}{\sum_k \sum_i q_{ik}} = \frac{\sum_i q_{ik}}{N}$ (erfordert einige Schritte mit Lagrange Multiplier und Ableitung)

EM für Gaußsche Mixture-Modelle

Gewichtete Maximum-Likelihood Updates:

- Koeffizienten aktualisieren: $\boldsymbol{\pi} = \arg \max_{\boldsymbol{\pi}} \sum_i \sum_k q_{ik} \log \pi_k$
 - Ergebnis: $\pi_k = \frac{\sum_i q_{ik}}{\sum_k \sum_i q_{ik}} = \frac{\sum_i q_{ik}}{N}$ (erfordert einige Schritte mit Lagrange Multiplier und Ableitung)
- Komponenten aktualisieren: $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k = \arg \max_{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k} \sum_i q_{ik} \log \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

EM für Gaußsche Mixture-Modelle

Gewichtete Maximum-Likelihood Updates:

- Koeffizienten aktualisieren: $\boldsymbol{\pi} = \arg \max_{\boldsymbol{\pi}} \sum_i \sum_k q_{ik} \log \pi_k$
 - Ergebnis: $\pi_k = \frac{\sum_i q_{ik}}{\sum_k \sum_i q_{ik}} = \frac{\sum_i q_{ik}}{N}$ (erfordert einige Schritte mit Lagrange Multiplier und Ableitung)
- Komponenten aktualisieren: $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k = \arg \max_{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k} \sum_i q_{ik} \log \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
 - Mittelwert: $\boldsymbol{\mu}_k = \frac{\sum_i q_{ik} \mathbf{x}_i}{\sum_i q_{ik}}$
 - Kovarianz: $\boldsymbol{\Sigma}_k = \frac{\sum_i q_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T}{\sum_i q_{ik}}$

EM für Gaußsche Mixture-Modelle

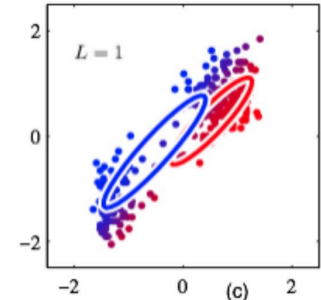
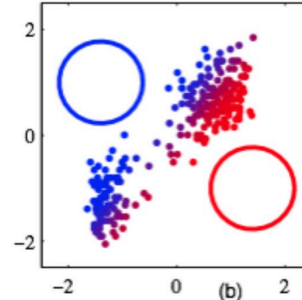
Gewichtete Maximum-Likelihood Updates:

- Koeffizienten aktualisieren: $\pi = \arg \max_{\pi} \sum_i \sum_k q_{ik} \log \pi_k$
 - Ergebnis: $\pi_k = \frac{\sum_i q_{ik}}{\sum_k \sum_i q_{ik}} = \frac{\sum_i q_{ik}}{N}$ (erfordert einige Schritte mit Lagrange Multiplier und Ableitung)

- Komponenten aktualisieren: $\mu_k, \Sigma_k = \arg \max_{\mu_k, \Sigma_k} \sum_i q_{ik} \log \mathcal{N}(\mathbf{x}_i | \mu_k, \Sigma_k)$

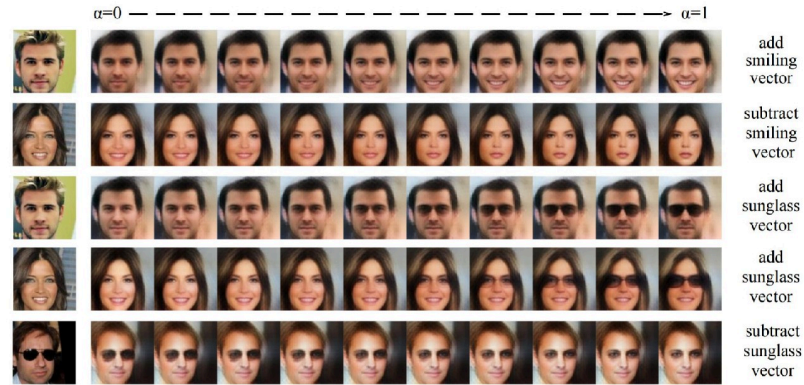
- Mittelwert:
$$\mu_k = \frac{\sum_i q_{ik} \mathbf{x}_i}{\sum_i q_{ik}}$$

- Kovarianz:
$$\Sigma_k = \frac{\sum_i q_{ik} (\mathbf{x}_i - \mu_k)(\mathbf{x}_i - \mu_k)^T}{\sum_i q_{ik}}$$



Dieselben Gleichungen wie zuvor, aber jetzt haben wir sie hergeleitet!

Teil 3



1) Expectation Maximization und Gaussian Mixture Models

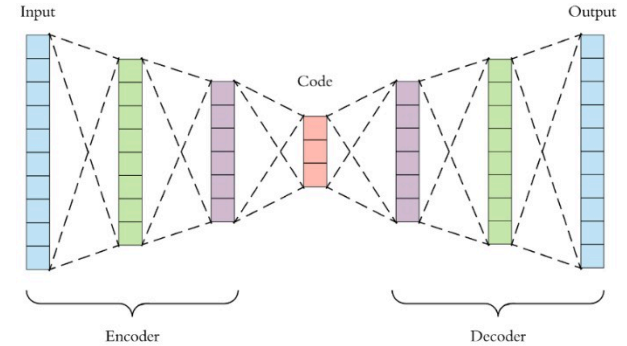
2) Verallgemeinerte Sicht auf EM

3) Variational Autoencoder

Variational-Autoencoder

Auto-Encoder

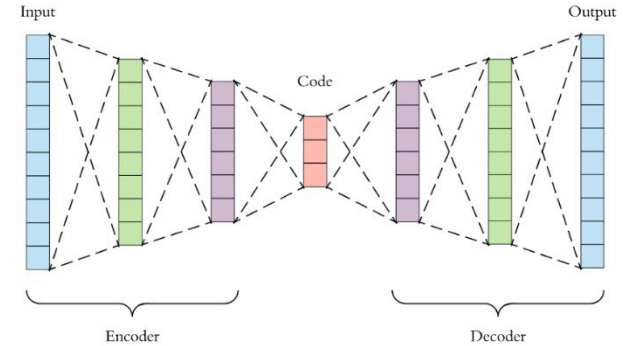
- Reduziert die Dimensionalität...
- ... aber es gibt keine Dichte
- ... und der AE kann keine neuen Daten generieren



Variational-Autoencoder

Auto-Encoder

- Reduziert die Dimensionalität...
- ... aber es gibt keine Dichte
- ... und der AE kann keine neuen Daten generieren

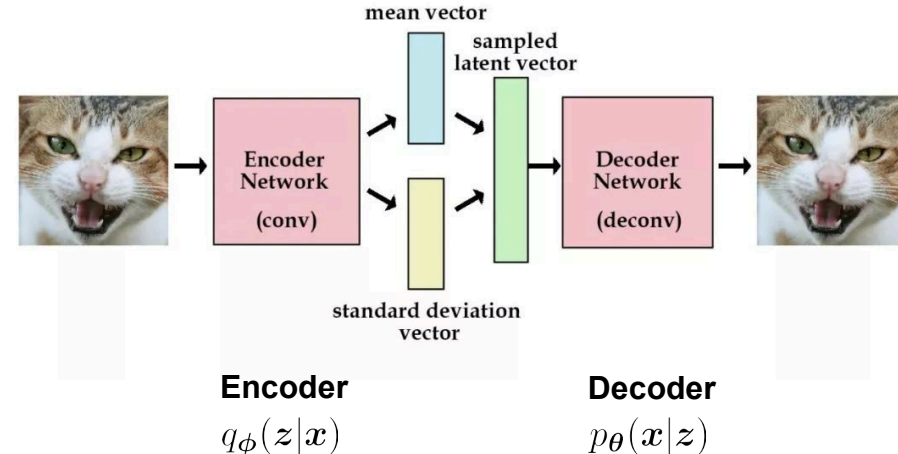


Variational-Autoencoder

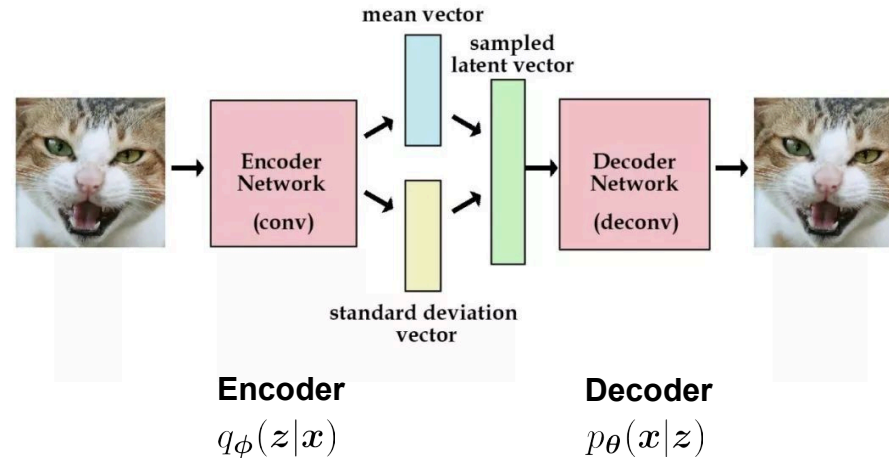
- Definiert auch eine Verteilung im latenten Raum
- **Generatives Modell:** Kann Samples im latenten Raum erzeugen
- ... und projizieren sie in den High-D-Raum (z. B. Bilder)
- Liefert eine **Dichte** (auch wenn sie sehr schwer zu berechnen ist)

Gelernt durch Variational Bayes:

- Verallgemeinerung von EM



Variational-Autoencoder: Einfache Erklärung



Zwei Loss Terme

- Reconstruction: Genau wie bei “normalem” Autoencoder: Vergleich von $p_{\theta}(x | z)$ und x
- Vergleich der latenten Verteilung $q_{\phi}(z|x)$ mit einer Zielverteilung, normalerweise $\mathcal{N}(0, 1)$

$q_{\phi}(z | x)$
↓
gesampled aus encoder

VAE Herleitung über amortisierte Variational Inference

Untere Schranke von EM:
$$\mathcal{L}(q, \theta) = \sum_{i=1}^N \mathbb{E}_{q_i(z_i)} [\log p_{\theta}(\mathbf{x}_i, z_i)] - \mathbb{E}_{q_i(z_i)} [\log q_i(z_i)]$$

VAE Herleitung über amortisierte Variational Inference

Untere Schranke von EM:
$$\mathcal{L}(q, \theta) = \sum_{i=1}^N \mathbb{E}_{q_i(z_i)} [\log p_{\theta}(\mathbf{x}_i, z_i)] - \mathbb{E}_{q_i(z_i)} [\log q_i(z_i)]$$

Anstatt eine individuelle Verteilung $q_i(z)$ pro Datenpunkt \mathbf{x}_i zu verwenden, können wir eine "amortisierte" Verteilung $q_{\phi}(z|\mathbf{x}_i)$ verwenden, die durch ein NN gegeben ist

VAE Herleitung über amortisierte Variational Inference

Untere Schranke von EM:
$$\mathcal{L}(q, \theta) = \sum_{i=1}^N \mathbb{E}_{q_i(z_i)} [\log p_{\theta}(\mathbf{x}_i, z_i)] - \mathbb{E}_{q_i(z_i)} [\log q_i(z_i)]$$

Anstatt eine individuelle Verteilung $q_i(z)$ pro Datenpunkt \mathbf{x}_i zu verwenden, können wir eine "amortisierte" Verteilung $q_{\phi}(z|\mathbf{x}_i)$ verwenden, die durch ein NN gegeben ist

$$\mathcal{L}(q, \theta) = \sum_{i=1}^N \mathbb{E}_{q_{\phi}(z|\mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i | z) p(z)] - \mathbb{E}_{q_{\phi}(z|\mathbf{x}_i)} [\log q_{\phi}(z | \mathbf{x}_i)]$$

VAE Herleitung über amortisierte Variational Inference

Untere Schranke von EM:
$$\mathcal{L}(q, \theta) = \sum_{i=1}^N \mathbb{E}_{q_i(z_i)} [\log p_{\theta}(\mathbf{x}_i, z_i)] - \mathbb{E}_{q_i(z_i)} [\log q_i(z_i)]$$

Anstatt eine individuelle Verteilung $q_i(z)$ pro Datenpunkt \mathbf{x}_i zu verwenden, können wir eine "amortisierte" Verteilung $q_{\phi}(z|\mathbf{x}_i)$ verwenden, die durch ein NN gegeben ist

$$\mathcal{L}(q, \theta) = \sum_{i=1}^N \mathbb{E}_{q_{\phi}(z|\mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i | z) p(z)] - \mathbb{E}_{q_{\phi}(z|\mathbf{x}_i)} [\log q_{\phi}(z | \mathbf{x}_i)]$$

$$\mathcal{L}(q, \theta) = \sum_{i=1}^N \mathbb{E}_{q_{\phi}(z|\mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i | z)] + \mathbb{E}_{q_{\phi}(z|\mathbf{x}_i)} [\log p(z)] - \mathbb{E}_{q_{\phi}(z|\mathbf{x}_i)} [\log q_{\phi}(z | \mathbf{x}_i)]$$

$$\mathcal{L}(q, \theta) = \sum_{i=1}^N \mathbb{E}_{q_{\phi}(z|\mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i | z)] + \mathbb{E}_{q_{\phi}(z|\mathbf{x}_i)} [\log p(z)] - \log q_{\phi}(z | \mathbf{x}_i)$$

Reconstruction loss KL-divergence loss

$\mathcal{N}(0, 1)$

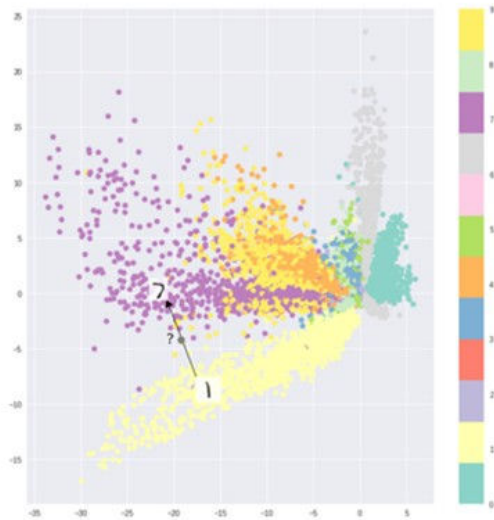
Dies ist das standard Objektiv für Variational Auto-Encoder (VAE)

- Encoder $q_{\phi}(z|\mathbf{x})$
- Decoder $p_{\theta}(\mathbf{x}|z)$

Visualisierung des Latenten Raums

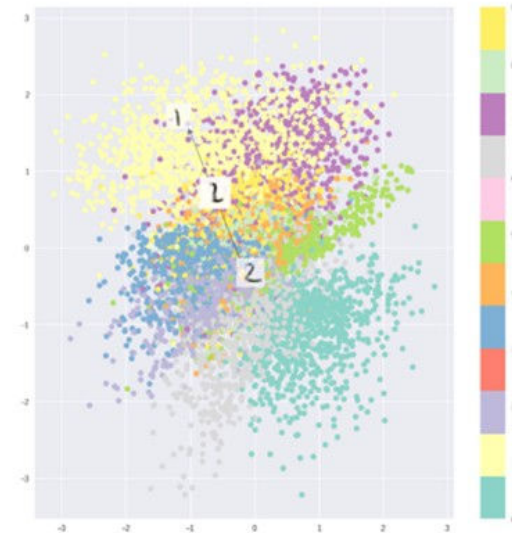
Auto-Encoder:

- Latenter Raum ist nicht zusammenhängend
- Keine sinnvolle Interpolation möglich



Variational Auto-Encoder:

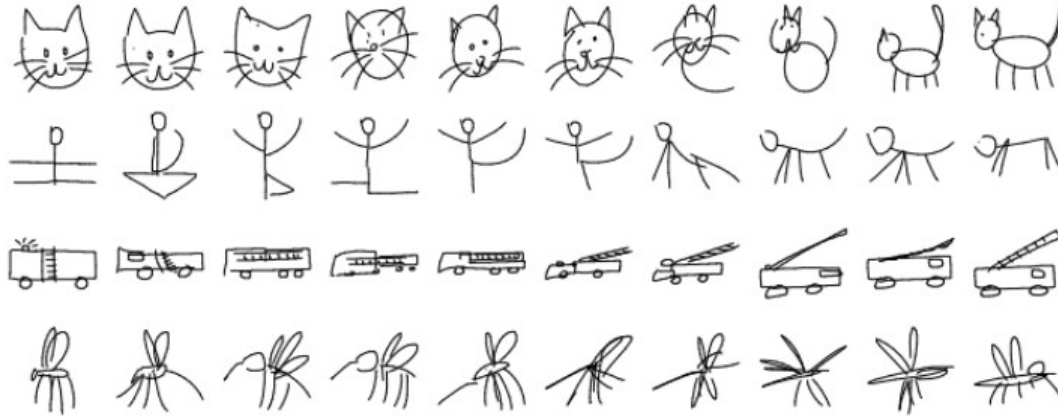
- Probabilistischer latenter Raum wird besser „gefüllt“
- Ermöglicht Interpolation



From: <https://www.jeremyjordan.me/variational-autoencoders/>

Interpolationen des latenten Raums

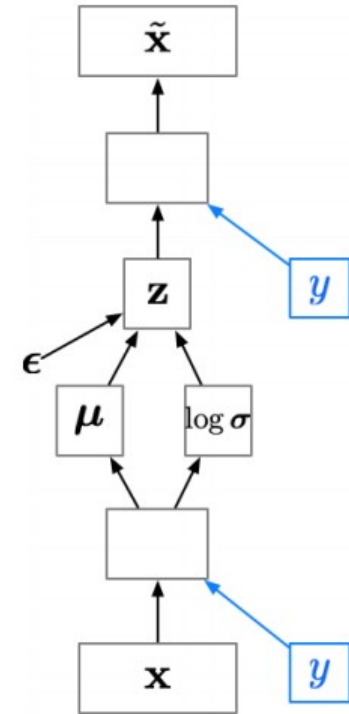
- Durch Interpolation zwischen zwei Vektoren im latenten Raum lassen sich oft interessante Ergebnisse erzielen:



Ha and Eck, "A neural representation of sketch drawings"

Klasse bedingte VAEs

- Bisher haben wir die Kennzeichnungen y nicht verwendet. Eine **klassenbedingte VAE** liefert die Kennzeichnungen sowohl an den Encoder als auch an den Decoder.
- Da der latente Code z nicht mehr die Bildkategorie modellieren muss, kann er sich auf die Modellierung der stilistischen Merkmale konzentrieren.
- Wenn wir Glück haben, können wir auf diese Weise Stil und Inhalt voneinander trennen. (Hinweis: Die Entflechtung ist immer noch eine dunkle Kunst.)
 - Siehe Kingma et al. "Semi-supervised learning with deep generative models".



Klassenbedingte VAE

- Indem man zwei **latente Dimensionen** (d.h. die Dimensionen von z) variiert, während man y konstant hält, kann man den **latenten Raum** visualisieren.



Klassenbedingte VAE

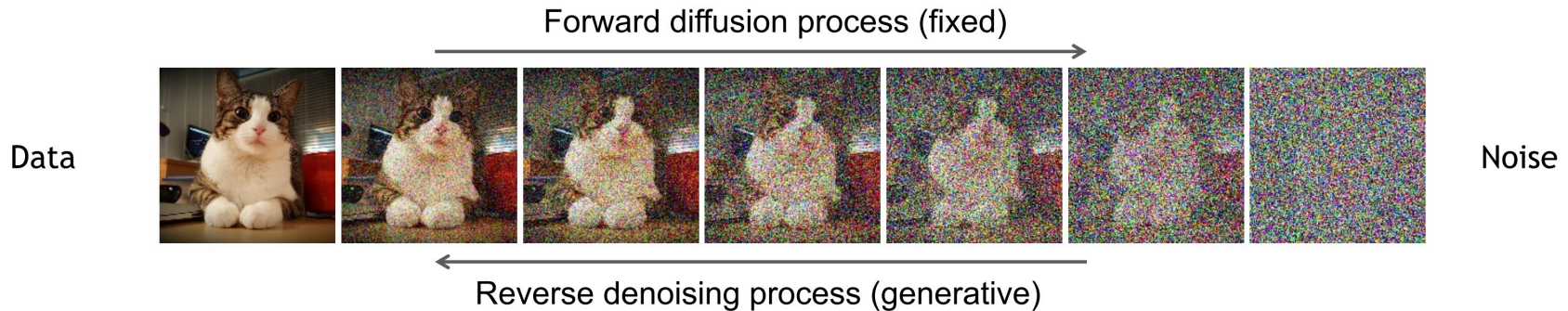
- Indem wir die Beschriftung y variieren, während wir z konstant halten, können wir Bildanalogien lösen



Outlook: Denoising Diffusion Probabilistic Models (DDPMs)

(Hierarchical) VAEs are the basis for diffusion models:

- Forward diffusion process that gradually adds noise to input
- Reverse denoising process that learns to generate data by denoising
- SOTA in image generation
- Covered in Machine Learning Lecture...



Fragen zum Selbsttest

- Was sind Mixture-Modelle?
- Sollten Gradientenmethoden für das Training von Mixture-Modellen verwendet werden?
- Wie funktioniert der EM-Algorithmus?
- Was optimiert der EM Algorithmus?
- Warum verbessert EM immer die untere Schranke?
- Warum können wir jede Komponente unabhängig voneinander mit EM optimieren?

