# Introduction to Generative Models
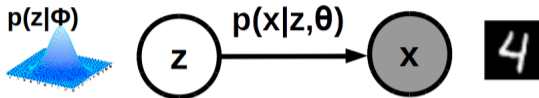
Piyush Rai

Machine Learning (CS771A)

Sept 23, 2016
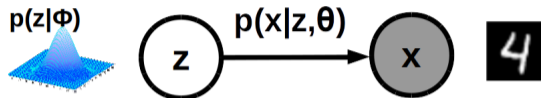
# Generative Model

- Defines a probabilistic way that could have "generated" the data
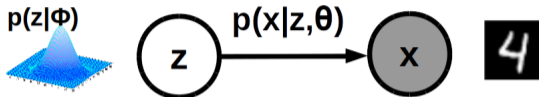
# Generative Model

- Defines a probabilistic way that could have "generated" the data



- Each observation $x_n$ is assumed to be associated with a latent variable $z_n$ (we can think of $z_n$ as a compact/compressed "encoding" of $x_n$)
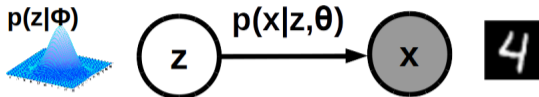
# Generative Model

- Defines a probabilistic way that could have "generated" the data



- Each observation $x_n$ is assumed to be associated with a latent variable $z_n$ (we can think of $z_n$ as a compact/compressed "encoding" of $x_n$)

- $z_n$ is assumed to be a random variable with some prior distribution $p(z|\phi)$

# Generative Model

- Defines a probabilistic way that could have "generated" the data



- Each observation $x_n$ is assumed to be associated with a latent variable $z_n$ (we can think of $z_n$ as a compact/compressed "encoding" of $x_n$)

- $z_n$ is assumed to be a random variable with some prior distribution $p(z|\phi)$

- Assume another data distribution $p(x|z, \theta)$ that can "generate" $x$ given $z$
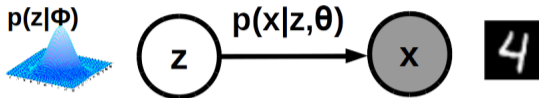
# Generative Model

- Defines a probabilistic way that could have "generated" the data



- Each observation $x_n$ is assumed to be associated with a latent variable $z_n$ (we can think of $z_n$ as a compact/compressed "encoding" of $x_n$)

- $z_n$ is assumed to be a random variable with some prior distribution $p(z|\phi)$

- Assume another data distribution $p(x|z, \theta)$ that can "generate" $x$ given $z$

- What $x$ and $z$ "look like", and the form of the distributions $p(z|\phi), p(x|z, \theta)$ will be problem-specific (we will soon look at some examples)

# Generative Model

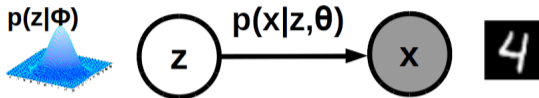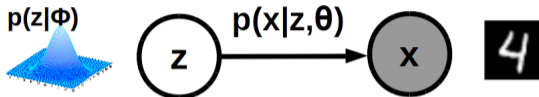- Defines a probabilistic way that could have "generated" the data



- Each observation $x_n$ is assumed to be associated with a latent variable $z_n$ (we can think of $z_n$ as a compact/compressed "encoding" of $x_n$)

- $z_n$ is assumed to be a random variable with some prior distribution $p(z|\phi)$

- Assume another data distribution $p(x|z, \theta)$ that can "generate" $x$ given $z$

- What $x$ and $z$ "look like", and the form of the distributions $p(z|\phi), p(x|z, \theta)$ will be problem-specific (we will soon look at some examples)

- $\{\theta, \phi\}$ are the unknown model parameters

# Generative Model

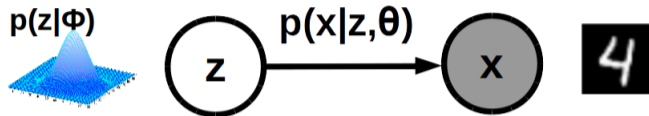- Defines a probabilistic way that could have "generated" the data



- Each observation $x_n$ is assumed to be associated with a latent variable $z_n$ (we can think of $z_n$ as a compact/compressed "encoding" of $x_n$)

- $z_n$ is assumed to be a random variable with some prior distribution $p(z|\phi)$

- Assume another data distribution $p(x|z, \theta)$ that can "generate" $x$ given $z$

- What $x$ and $z$ "look like", and the form of the distributions $p(z|\phi), p(x|z, \theta)$ will be problem-specific (we will soon look at some examples)

- $\{\theta, \phi\}$ are the unknown model parameters

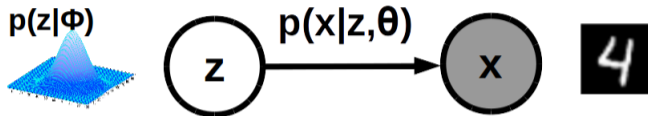- The goal will be to learn $\{\theta, \phi\}$ and $z_n$'s, given the observed data

# "Generative Story" of Data

- Generative models can be described using a "generative story" for the data

# "Generative Story" of Data

- Generative models can be described using a "generative story" for the data



$$p(z|\Phi) \qquad z \qquad p(x|z,\theta) \qquad x$$

- The "generative story" of each observation $\boldsymbol{x}_n$, $\forall n$

# "Generative Story" of Data

- Generative models can be described using a "generative story" for the data



$$p(z|\Phi) \qquad z \qquad p(x|z,\theta) \qquad x$$

- The "generative story" of each observation $x_n$, $\forall n$
  - First draw a random latent variable $z_n \sim p(z|\phi)$ from the prior on $z$

# "Generative Story" of Data

- Generative models can be described using a "generative story" for the data



$$p(z|\Phi) \qquad \underset{z}{\bigcirc} \quad \xrightarrow{\ p(x|z,\theta)\ } \quad \underset{x}{\bigcirc}$$

- The "generative story" of each observation $x_n$, $\forall n$
  - First draw a random latent variable $z_n \sim p(z|\phi)$ from the prior on $z$
  - Given $z_n$, now generate $x_n$ as $x_n \sim p(x|\theta, z_n)$ from the data distribution
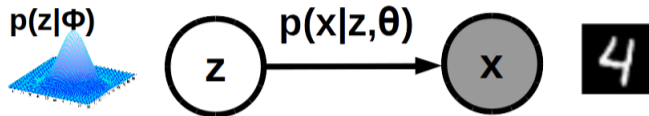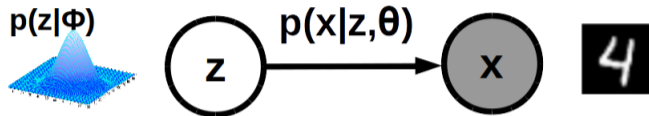
# "Generative Story" of Data

- Generative models can be described using a "generative story" for the data



- The "generative story" of each observation $x_n$, $\forall n$
  - First draw a random latent variable $z_n \sim p(z|\phi)$ from the prior on $z$
  - Given $z_n$, now generate $x_n$ as $x_n \sim p(x|\theta, z_n)$ from the data distribution
- Such models usually have two types of variables: "local" and "global"
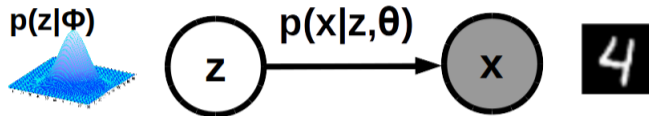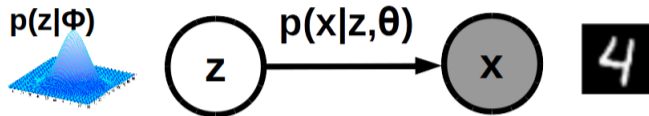
# "Generative Story" of Data

- Generative models can be described using a "generative story" for the data



- The "generative story" of each observation $x_n$, $\forall n$
  - First draw a random latent variable $z_n \sim p(z|\phi)$ from the prior on $z$
  - Given $z_n$, now generate $x_n$ as $x_n \sim p(x|\theta, z_n)$ from the data distribution
- Such models usually have two types of variables: "local" and "global"
  - Each $z_n$ is a "local" variable (specific to the data point $x_n$)
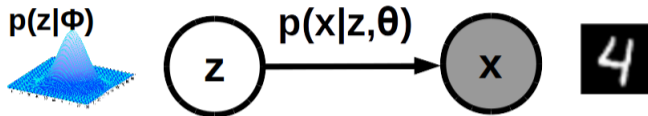
# "Generative Story" of Data

- Generative models can be described using a "generative story" for the data



- The "generative story" of each observation $x_n$, $\forall n$
  - First draw a random latent variable $z_n \sim p(z|\phi)$ from the prior on $z$
  - Given $z_n$, now generate $x_n$ as $x_n \sim p(x|\theta, z_n)$ from the data distribution
- Such models usually have two types of variables: "local" and "global"
  - Each $z_n$ is a "local" variable (specific to the data point $x_n$)
  - $(\phi, \theta)$ are global variables (shared by all the data points)
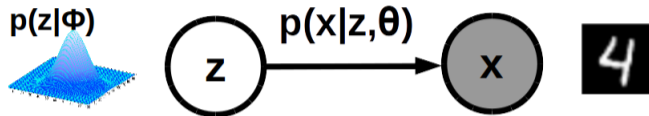
# "Generative Story" of Data

- Generative models can be described using a "generative story" for the data



- The "generative story" of each observation $x_n$, $\forall n$
  - First draw a random latent variable $z_n \sim p(z|\phi)$ from the prior on $z$
  - Given $z_n$, now generate $x_n$ as $x_n \sim p(x|\theta, z_n)$ from the data distribution
- Such models usually have two types of variables: "local" and "global"
  - Each $z_n$ is a "local" variable (specific to the data point $x_n$)
  - $(\phi, \theta)$ are global variables (shared by all the data points)
  - We may be interested in learning the global vars, or local vars, or both

# "Generative Story" of Data

- Generative models can be described using a "generative story" for the data



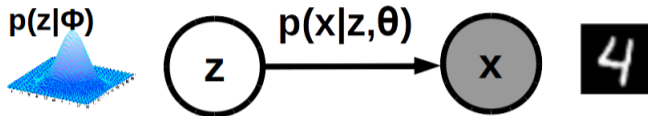- The "generative story" of each observation $x_n$, $\forall n$
    - First draw a random latent variable $z_n \sim p(z|\phi)$ from the prior on $z$
    - Given $z_n$, now generate $x_n$ as $x_n \sim p(x|\theta, z_n)$ from the data distribution
- Such models usually have two types of variables: "local" and "global"
    - Each $z_n$ is a "local" variable (specific to the data point $x_n$)
    - $(\phi, \theta)$ are global variables (shared by all the data points)
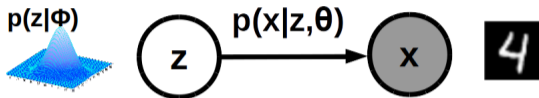    - We may be interested in learning the global vars, or local vars, or both
    - Usually it's possible to infer the global vars from local vars (or vice-versa)
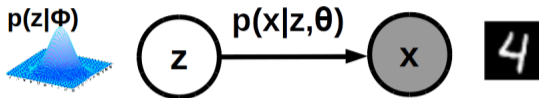
# Why Generative Models?

- A proper, probabilistic way to think about the data generation process

# Why Generative Models?

- A proper, probabilistic way to think about the data generation process



- Allows modeling different types of data (real, binary, count, etc.) by changing the data distribution $p(\boldsymbol{x}|\theta, \boldsymbol{z})$ appropriately

# Why Generative Models?

- A proper, probabilistic way to think about the data generation process



- Allows modeling different types of data (real, binary, count, etc.) by changing the data distribution $p(\boldsymbol{x}|\theta, \boldsymbol{z})$ appropriately

- Can synthesize or "hallucinate" new data using an already learned model

  - Generate a "random" $\boldsymbol{z}$ from $p(\boldsymbol{z}|\phi)$ and generate $\boldsymbol{x}$ from $p(\boldsymbol{x}|\theta, \boldsymbol{z})$
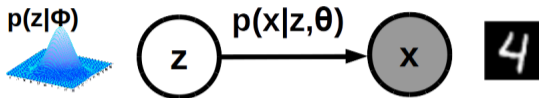
# Why Generative Models?

- A proper, probabilistic way to think about the data generation process



- Allows modeling different types of data (real, binary, count, etc.) by changing the data distribution $p(\boldsymbol{x}|\theta, \boldsymbol{z})$ appropriately

- Can synthesize or "hallucinate" new data using an already learned model
  - Generate a "random" $\boldsymbol{z}$ from $p(\boldsymbol{z}|\phi)$ and generate $\boldsymbol{x}$ from $p(\boldsymbol{x}|\theta, \boldsymbol{z})$



after 5 epochs      after 100 epochs

- Allows handling missing data (by treating missing data also as latent variable)

# Some "Canonical" Generative Models

- Mixture model (used in clustering and probability density estimation)

# Some "Canonical" Generative Models

- Mixture model (used in clustering and probability density estimation)



- Latent factor model (used in dimensionality reduction)
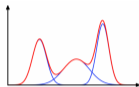
# Some "Canonical" Generative Models

- Mixture model (used in clustering and probability density estimation)



- Latent factor model (used in dimensionality reduction)



- Can even combine these (e.g., mixture of latent factor models)

# Example: Mixture Model

- Assume data $\{\boldsymbol{x}_n\}_{n=1}^{N}$ was generated from a mixture of $K$ distributions

## Example: Mixture Model

- Assume data $\{\boldsymbol{x}_n\}_{n=1}^{N}$ was generated from a mixture of $K$ distributions



- Suppose these $K$ distributions are $p(\boldsymbol{x}|\theta_1), \ldots, p(\boldsymbol{x}|\theta_K)$

# Example: Mixture Model

- Assume data $\{x_n\}_{n=1}^{N}$ was generated from a mixture of $K$ distributions



- Suppose these $K$ distributions are $p(x|\theta_1), \ldots, p(x|\theta_K)$
- Don't know which of the $K$ distributions each $x_n$ was generated from

## Example: Mixture Model

- Assume data $\{\boldsymbol{x}_n\}_{n=1}^N$ was generated from a mixture of $K$ distributions



- Suppose these $K$ distributions are $p(\boldsymbol{x}|\theta_1), \ldots, p(\boldsymbol{x}|\theta_K)$
- Don't know which of the $K$ distributions each $\boldsymbol{x}_n$ was generated from
- Consider the following generative story for each $\boldsymbol{x}_n$, $n = 1, 2, \ldots, N$

# Example: Mixture Model

- Assume data $\{\boldsymbol{x}_n\}_{n=1}^N$ was generated from a mixture of $K$ distributions



- Suppose these $K$ distributions are $p(\boldsymbol{x}|\theta_1), \ldots, p(\boldsymbol{x}|\theta_K)$
- Don't know which of the $K$ distributions each $\boldsymbol{x}_n$ was generated from
- Consider the following generative story for each $\boldsymbol{x}_n$, $n = 1, 2, \ldots, N$
  - First choose a mixture component $\boldsymbol{z}_n \in \{1, 2, \ldots, K\}$ as $\boldsymbol{z}_n \sim p(\boldsymbol{z}|\phi)$

# Example: Mixture Model
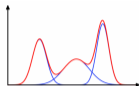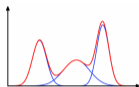
- Assume data $\{x_n\}_{n=1}^{N}$ was generated from a mixture of $K$ distributions



- Suppose these $K$ distributions are $p(x|\theta_1), \ldots, p(x|\theta_K)$

- Don't know which of the $K$ distributions each $x_n$ was generated from

- Consider the following generative story for each $x_n$, $n = 1, 2, \ldots, N$
  - First choose a mixture component $z_n \in \{1, 2, \ldots, K\}$ as $z_n \sim p(z|\phi)$
  - Now generate $x_n$ from the mixture component no. $z_n$ as $x_n \sim p(x|\theta_{z_n})$

## Example: Mixture Model
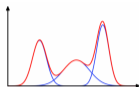
- Assume data $\{x_n\}_{n=1}^N$ was generated from a mixture of $K$ distributions



- Suppose these $K$ distributions are $p(x|\theta_1), \ldots, p(x|\theta_K)$

- Don't know which of the $K$ distributions each $x_n$ was generated from

- Consider the following generative story for each $x_n$, $n = 1, 2, \ldots, N$
  - First choose a mixture component $z_n \in \{1, 2, \ldots, K\}$ as $z_n \sim p(z|\phi)$
  - Now generate $x_n$ from the mixture component no. $z_n$ as $x_n \sim p(x|\theta_{z_n})$

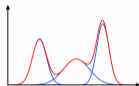- In a mixture model, $z$ is discrete so $p(z|\phi)$ is a multinomial distribution

## Example: Mixture Model

- Assume data $\{x_n\}_{n=1}^{N}$ was generated from a mixture of $K$ distributions



- Suppose these $K$ distributions are $p(x|\theta_1), \ldots, p(x|\theta_K)$

- Don't know which of the $K$ distributions each $x_n$ was generated from

- Consider the following generative story for each $x_n$, $n = 1, 2, \ldots, N$

  - First choose a mixture component $z_n \in \{1, 2, \ldots, K\}$ as $z_n \sim p(z|\phi)$

  - Now generate $x_n$ from the mixture component no. $z_n$ as $x_n \sim p(x|\theta_{z_n})$

- In a mixture model, $z$ is discrete so $p(z|\phi)$ is a multinomial distribution

- The data distribution $p(x|\theta_{z_n})$ depends on the type of data being modeled
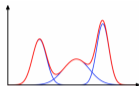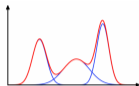
# Example: Mixture Model

- Assume data $\{\boldsymbol{x}_n\}_{n=1}^N$ was generated from a mixture of $K$ distributions
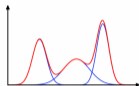


- Suppose these $K$ distributions are $p(\boldsymbol{x}|\theta_1), \ldots, p(\boldsymbol{x}|\theta_K)$

- Don't know which of the $K$ distributions each $\boldsymbol{x}_n$ was generated from

- Consider the following generative story for each $\boldsymbol{x}_n$, $n = 1, 2, \ldots, N$
  - First choose a mixture component $\boldsymbol{z}_n \in \{1, 2, \ldots, K\}$ as $\boldsymbol{z}_n \sim p(\boldsymbol{z}|\phi)$
  - Now generate $\boldsymbol{x}_n$ from the mixture component no. $\boldsymbol{z}_n$ as $\boldsymbol{x}_n \sim p(\boldsymbol{x}|\theta_{\boldsymbol{z}_n})$

- In a mixture model, $\boldsymbol{z}$ is discrete so $p(\boldsymbol{z}|\phi)$ is a multinomial distribution

- The data distribution $p(\boldsymbol{x}|\theta_{\boldsymbol{z}_n})$ depends on the type of data being modeled

- Mixture models can model complex distributions as superposition of simpler distributions (can be used for density estimation, as well as clustering).

# Example: Latent Factor Model

- Assume each $D$-dim $\boldsymbol{x}_n$ generated from a $K$-dim latent factor $\boldsymbol{z}_n$ ($K \ll D$)

## Example: Latent Factor Model

- Assume each $D$-dim $\boldsymbol{x}_n$ generated from a $K$-dim latent factor $\boldsymbol{z}_n$ ($K \ll D$)



- Consider the following generative story for each $\boldsymbol{x}_n$, $n = 1, 2, \ldots, N$

## Example: Latent Factor Model

- Assume each $D$-dim $\boldsymbol{x}_n$ generated from a $K$-dim latent factor $\boldsymbol{z}_n$ ($K \ll D$)



- Consider the following generative story for each $\boldsymbol{x}_n$, $n = 1, 2, \ldots, N$
  - First generate $\boldsymbol{z}_n$ from a $K$-dim distr. as $\boldsymbol{z}_n \sim p(\boldsymbol{z}|\phi)$

## Example: Latent Factor Model

- Assume each $D$-dim $\boldsymbol{x}_n$ generated from a $K$-dim latent factor $\boldsymbol{z}_n$ ($K \ll D$)



- Consider the following generative story for each $\boldsymbol{x}_n$, $n = 1, 2, \ldots, N$
  - First generate $\boldsymbol{z}_n$ from a $K$-dim distr. as $\boldsymbol{z}_n \sim p(\boldsymbol{z}|\phi)$
  - Now generate $\boldsymbol{x}_n$ from a $D$-dim distr. as $\boldsymbol{x}_n \sim p(\boldsymbol{x}|\boldsymbol{z}_n, \theta)$
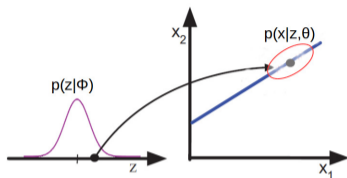
# Example: Latent Factor Model

- Assume each $D$-dim $\boldsymbol{x}_n$ generated from a $K$-dim latent factor $\boldsymbol{z}_n$ ($K \ll D$)



- Consider the following generative story for each $\boldsymbol{x}_n$, $n = 1, 2, \ldots, N$
  - First generate $\boldsymbol{z}_n$ from a $K$-dim distr. as $\boldsymbol{z}_n \sim p(\boldsymbol{z}|\phi)$
  - Now generate $\boldsymbol{x}_n$ from a $D$-dim distr. as $\boldsymbol{x}_n \sim p(\boldsymbol{x}|\boldsymbol{z}_n, \theta)$
- When $p(\boldsymbol{z}|\phi)$ and $p(\boldsymbol{x}|\boldsymbol{z}_n, \theta)$ are Gaussian distributions, this basic generative model is called factor analysis or probabilistic PCA
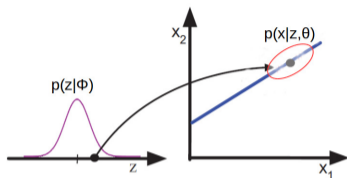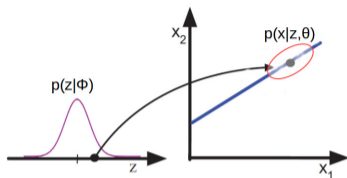
# Example: Latent Factor Model

- Assume each $D$-dim $\boldsymbol{x}_n$ generated from a $K$-dim latent factor $\boldsymbol{z}_n$ ($K \ll D$)



- Consider the following generative story for each $\boldsymbol{x}_n$, $n = 1, 2, \ldots, N$
  - First generate $\boldsymbol{z}_n$ from a $K$-dim distr. as $\boldsymbol{z}_n \sim p(\boldsymbol{z}|\phi)$
  - Now generate $\boldsymbol{x}_n$ from a $D$-dim distr. as $\boldsymbol{x}_n \sim p(\boldsymbol{x}|\boldsymbol{z}_n, \theta)$

- When $p(\boldsymbol{z}|\phi)$ and $p(\boldsymbol{x}|\boldsymbol{z}_n, \theta)$ are Gaussian distributions, this basic generative model is called factor analysis or probabilistic PCA

- The choice of $p(\boldsymbol{z}|\phi)$ and $p(\boldsymbol{x}|\boldsymbol{z}_n, \theta)$ in general will be problem dependent
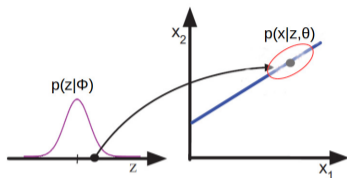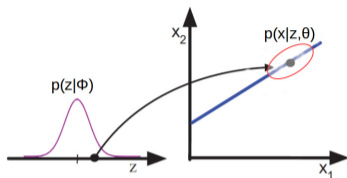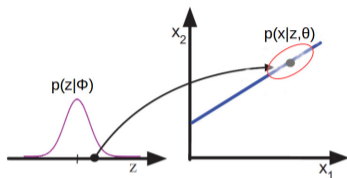
# Example: Latent Factor Model

- Assume each $D$-dim $\boldsymbol{x}_n$ generated from a $K$-dim latent factor $\boldsymbol{z}_n$ ($K \ll D$)



- Consider the following generative story for each $\boldsymbol{x}_n$, $n = 1, 2, \ldots, N$
  - First generate $\boldsymbol{z}_n$ from a $K$-dim distr. as $\boldsymbol{z}_n \sim p(\boldsymbol{z}|\phi)$
  - Now generate $\boldsymbol{x}_n$ from a $D$-dim distr. as $\boldsymbol{x}_n \sim p(\boldsymbol{x}|\boldsymbol{z}_n, \theta)$
- When $p(\boldsymbol{z}|\phi)$ and $p(\boldsymbol{x}|\boldsymbol{z}_n, \theta)$ are Gaussian distributions, this basic generative model is called factor analysis or probabilistic PCA
- The choice of $p(\boldsymbol{z}|\phi)$ and $p(\boldsymbol{x}|\boldsymbol{z}_n, \theta)$ in general will be problem dependent
- Many recent advances in generative models (e.g., deep generative models, generative adversarial networks, etc) are based on these basic principles

## Going Forward..

- We will look at, in more detail, some specific generative models

# Going Forward..

- We will look at, in more detail, some specific generative models
  - Gaussian mixture model (for clustering and density estimation)

## Going Forward..

- We will look at, in more detail, some specific generative models
  - Gaussian mixture model (for clustering and density estimation)
  - Factor Analysis and Probabilistic PCA (for dimensionality reduction)

## Going Forward..

- We will look at, in more detail, some specific generative models
  - Gaussian mixture model (for clustering and density estimation)
  - Factor Analysis and Probabilistic PCA (for dimensionality reduction)
- We will also look at how to do parameter estimation in such models

## Going Forward..

- We will look at, in more detail, some specific generative models
  - Gaussian mixture model (for clustering and density estimation)
  - Factor Analysis and Probabilistic PCA (for dimensionality reduction)
- We will also look at how to do parameter estimation in such models
  - One common approach is to perform MLE/MAP

# Going Forward..

- We will look at, in more detail, some specific generative models
  - Gaussian mixture model (for clustering and density estimation)
  - Factor Analysis and Probabilistic PCA (for dimensionality reduction)
- We will also look at how to do parameter estimation in such models
  - One common approach is to perform MLE/MAP
  - However, presence of latent variables $z$ makes MLE/MAP hard

# Going Forward..

- We will look at, in more detail, some specific generative models
    - Gaussian mixture model (for clustering and density estimation)
    - Factor Analysis and Probabilistic PCA (for dimensionality reduction)
- We will also look at how to do parameter estimation in such models
    - One common approach is to perform MLE/MAP
    - However, presence of latent variables $z$ makes MLE/MAP hard
    - Reason: Since $z$ is a random variable, we must sum over all possible values of $z$ when doing MLE/MAP for the model parameters $\theta, \phi$

$$\log p(\boldsymbol{x}|\theta, \phi) = \log \sum_{\boldsymbol{z}} p(\boldsymbol{x}|\boldsymbol{z}, \theta) p(\boldsymbol{z}|\phi) \quad \text{(Log can't go inside the summation!)}$$

# Going Forward..

- We will look at, in more detail, some specific generative models
  - Gaussian mixture model (for clustering and density estimation)
  - Factor Analysis and Probabilistic PCA (for dimensionality reduction)
- We will also look at how to do parameter estimation in such models
  - One common approach is to perform MLE/MAP
  - However, presence of latent variables $z$ makes MLE/MAP hard
  - Reason: Since $z$ is a random variable, we must sum over all possible values of $z$ when doing MLE/MAP for the model parameters $\theta, \phi$

  $$\log p(x|\theta, \phi) = \log \sum_z p(x|z, \theta) p(z|\phi) \quad \text{(Log can't go inside the summation!)}$$

  - Expectation Maximization (EM) algorithm gives a way to solve the problem
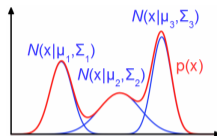
# Going Forward..

- We will look at, in more detail, some specific generative models
  - Gaussian mixture model (for clustering and density estimation)
  - Factor Analysis and Probabilistic PCA (for dimensionality reduction)
- We will also look at how to do parameter estimation in such models
  - One common approach is to perform MLE/MAP
  - However, presence of latent variables $z$ makes MLE/MAP hard
  - Reason: Since $z$ is a random variable, we must sum over all possible values of $z$ when doing MLE/MAP for the model parameters $\theta, \phi$

  $$\log p(\mathbf{x}|\theta, \phi) = \log \sum_{z} p(\mathbf{x}|\mathbf{z}, \theta) p(\mathbf{z}|\phi) \quad \text{(Log can't go inside the summation!)}$$

  - Expectation Maximization (EM) algorithm gives a way to solve the problem
  - Basic idea in EM: Instead of summing over all possibilities of $z$, make a "guess" $\tilde{z}$ and maximize $\log p(\mathbf{x}, \tilde{\mathbf{z}}|\theta, \phi)$ w.r.t. $\theta, \phi$ to learn $\theta, \phi$. Use these values of $\theta, \phi$ to refine your guess $\tilde{z}$ and repeat until convergence.
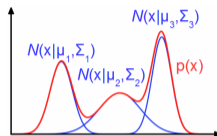
# Gaussian Mixture Model (GMM)

- Assume the data is generated from a mixture of $K$ Gaussians



- Each Gaussian represents a "cluster" in the data

# Gaussian Mixture Model (GMM)

- Assume the data is generated from a mixture of $K$ Gaussians



- Each Gaussian represents a "cluster" in the data

- The distribution $p(\boldsymbol{x})$ will be a weighted a mixture of $K$ Gaussians

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

# Gaussian Mixture Model (GMM)

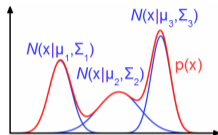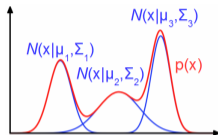- Assume the data is generated from a mixture of $K$ Gaussians



- Each Gaussian represents a "cluster" in the data
- The distribution $p(\boldsymbol{x})$ will be a weighted a mixture of $K$ Gaussians

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

where $\pi_k$'s are the **mixing weights**: $\sum_{k=1}^{K} \pi_k = 1, \pi_k \geq 0$ (intuitively, $\pi_k = p(\boldsymbol{z} = k)$ is the fraction of data generated by the $k$-th distribution)

# Gaussian Mixture Model (GMM)

- Assume the data is generated from a mixture of $K$ Gaussians



- Each Gaussian represents a "cluster" in the data

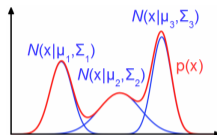- The distribution $p(\boldsymbol{x})$ will be a weighted a mixture of $K$ Gaussians

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

where $\pi_k$'s are the **mixing weights**: $\sum_{k=1}^{K} \pi_k = 1, \pi_k \geq 0$ (intuitively, $\pi_k = p(\boldsymbol{z} = k)$ is the fraction of data generated by the $k$-th distribution)

- The goal is to learn the params $\{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^{K}$ of these $K$ Gaussians, the mixing weights $\{\pi_k\}_{k=1}^{K}$, and/or the cluster assignment $\boldsymbol{z}_n$ of each $\boldsymbol{x}_n$

# Gaussian Mixture Model (GMM)

- Assume the data is generated from a mixture of $K$ Gaussians



- Each Gaussian represents a "cluster" in the data
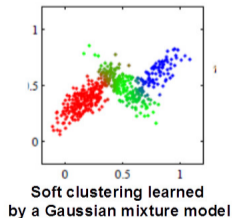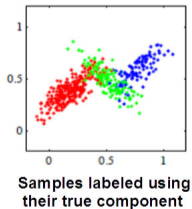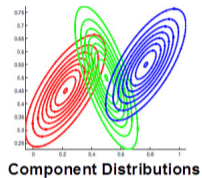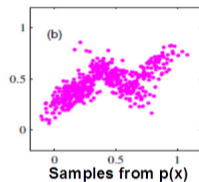- The distribution $p(\boldsymbol{x})$ will be a weighted a mixture of $K$ Gaussians

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

where $\pi_k$'s are the **mixing weights**: $\sum_{k=1}^{K} \pi_k = 1, \pi_k \geq 0$ (intuitively, $\pi_k = p(\boldsymbol{z} = k)$ is the fraction of data generated by the $k$-th distribution)

- The goal is to learn the params $\{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^{K}$ of these $K$ Gaussians, the mixing weights $\{\pi_k\}_{k=1}^{K}$, and/or the cluster assignment $\boldsymbol{z}_n$ of each $\boldsymbol{x}_n$
- GMM in many ways improves over $K$-means clustering

# GMM Clustering: Pictorially

Some synthetically generated data (top-left) generated from a mixture of 3 overlapping Gaussians (top-right).



Samples from p(x)

Component Distributions

Samples labeled using their true component

Soft clustering learned by a Gaussian mixture model

Notice the "mixed" colored points in the overlapping regions in the final clustering

## Next Class

- GMM in more detail. Extensions of GMM.

- Parameter estimation in GMM

- The Expectation Maximization (EM) algorithm