

Linear Factor Models

Linear Factor Models

- Defined by the use of a stochastic, linear decoder that generates \mathbf{x} by adding noise to a linear transformation of \mathbf{z} .
- LFMs describe the data generation process as follows:

Sample the explanatory factors \mathbf{z} from a distribution

$$\mathbf{z} \sim \mathbf{p}(\mathbf{z})$$

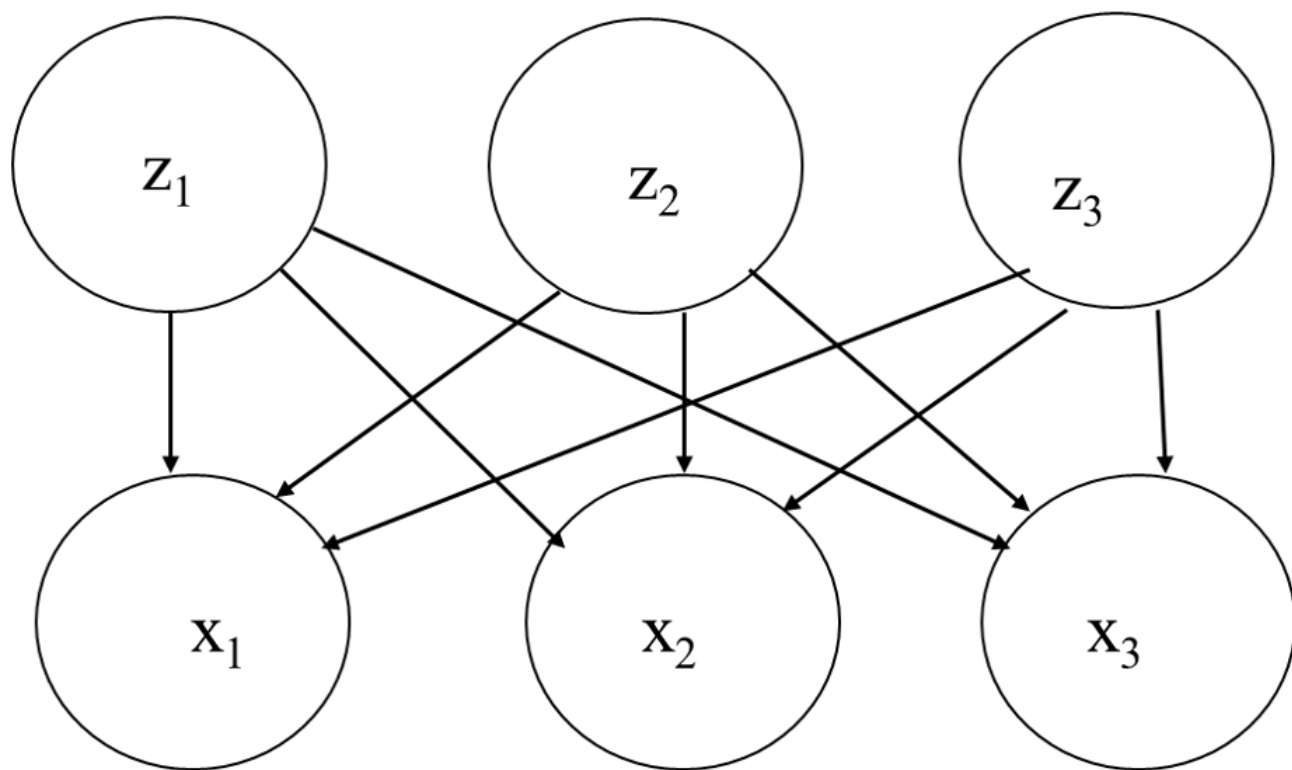
where $\mathbf{p}(\mathbf{z})$ is a factorial distribution (i.e. $\mathbf{p}(\mathbf{z}) = \prod_i \mathbf{p}(z_i)$)

- Sample the real-valued observable variables given the factors:

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \mathbf{noise}$$

where the noise is typically Gaussian and diagonal

Linear Factor Models



$$x = Wz + \mu + noise$$

Types of LFM

- The directed graphical model on the previous slide describes the LFM family, where we assume that observed \mathbf{x} is obtained by a linear combination of independent latent factors \mathbf{z} , plus some noise.
- Different types of LFM make different choices about the form of the noise and of the prior $p(\mathbf{z})$.
- We will touch upon:
 - a) Probabilistic PCA
 - b) Independent component analysis (ICA)
 - c) Slow feature analysis
 - d) Sparse coding

Probabilistic PCA

- Probabilistic PCA (Principal Component Analysis) is a probabilistic extension of the classical PCA method, which is widely used for dimensionality reduction and data visualization. PCA is typically used to find the orthogonal axes (principal components) along which the data varies the most. It projects the data onto these axes, effectively reducing the dimensionality of the data while retaining most of its variance.
- In probabilistic PCA, the model assumes that the observed data points are generated from a latent variable model with Gaussian noise. This latent variable model represents the underlying structure of the data. The main idea is to introduce latent variables that capture the variability of the data in a lower-dimensional space, and then model the observed data as noisy projections of these latent variables onto a higher-dimensional space.
- Probabilistic PCA is useful in scenarios where there is uncertainty in the data or when there is missing information. It provides a more flexible framework compared to classical PCA and can handle noise and missing data more effectively.

Probabilistic PCA

1. Model Setup:

We start with a dataset \mathbf{X} consisting of N observations of D -dimensional data points. Each data point can be represented as a vector \mathbf{x}_n , where $n = 1, 2, \dots, N$.

The generative model for PPCA assumes that the observed data points are generated from a lower-dimensional latent space. Mathematically, this is represented as:

$$\mathbf{x}_n = \mathbf{W}\mathbf{z}_n + \boldsymbol{\mu} + \boldsymbol{\epsilon}_n$$

where:

- \mathbf{W} is a $D \times L$ matrix of principal components, where L is the dimensionality of the latent space.
- \mathbf{z}_n is a L -dimensional latent variable vector corresponding to the n -th observation.
- $\boldsymbol{\mu}$ is the mean vector.
- $\boldsymbol{\epsilon}_n$ is a D -dimensional Gaussian noise vector with zero mean and covariance matrix $\sigma^2 \mathbf{I}$, where σ^2 is the noise variance and \mathbf{I} is the identity matrix.

Probabilistic PCA

2. Likelihood Function:

Given the generative model, the likelihood function of the observed data given the latent variables and parameters is assumed to be Gaussian:

$$p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{W}, \boldsymbol{\mu}, \sigma^2) = \mathcal{N}(\mathbf{x}_n | \mathbf{W}\mathbf{z}_n + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$$

where $\mathcal{N}(\mathbf{x}_n | \mathbf{W}\mathbf{z}_n + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$ represents the multivariate Gaussian distribution.

3. Latent Variable Distribution:

The latent variables \mathbf{z}_n are assumed to follow a Gaussian distribution:

$$p(\mathbf{z}_n) = \mathcal{N}(\mathbf{z}_n | \mathbf{0}, \mathbf{I})$$

where $\mathbf{0}$ represents a zero vector and \mathbf{I} represents the identity matrix.

4. Parameter Estimation:

Parameter estimation in PPCA typically involves maximizing the log-likelihood function with respect to the parameters \mathbf{W} , $\boldsymbol{\mu}$, and σ^2 .

Probabilistic PCA

5. Expectation-Maximization (EM) Algorithm:

- The EM algorithm is an iterative optimization technique used to estimate the parameters of probabilistic models when there are latent variables. In the case of PPCA, the EM algorithm alternates between two steps: the E-step and the M-step.
- E-step (Expectation Step):

In the E-step, we compute the expected values of the latent variables given the observed data and the current parameter estimates. Specifically, we compute the posterior distribution of the latent variables conditioned on the observed data and the current parameter estimates.

Probabilistic PCA

Given the observed data \mathbf{X} , the latent variable distribution $p(\mathbf{z}_n)$, and the current parameter estimates \mathbf{W} , $\boldsymbol{\mu}$, and σ^2 , the posterior distribution of the latent variables \mathbf{z}_n can be computed using Bayes' theorem:

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}, \boldsymbol{\mu}, \sigma^2) = \frac{p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{W}, \boldsymbol{\mu}, \sigma^2) \cdot p(\mathbf{z}_n)}{p(\mathbf{x}_n | \mathbf{W}, \boldsymbol{\mu}, \sigma^2)}$$

Where:

- $p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{W}, \boldsymbol{\mu}, \sigma^2)$ is the likelihood function of observing \mathbf{x}_n given \mathbf{z}_n and the parameters.
- $p(\mathbf{z}_n)$ is the prior distribution of the latent variables.
- $p(\mathbf{x}_n | \mathbf{W}, \boldsymbol{\mu}, \sigma^2)$ is the marginal likelihood of observing \mathbf{x}_n given the parameters, obtained by integrating over all possible values of \mathbf{z}_n .

Typically, in PPCA, the latent variable distribution $p(\mathbf{z}_n)$ is assumed to be Gaussian with zero mean and identity covariance matrix, as mentioned earlier.

After computing the posterior distribution of the latent variables, we use the expected values of the latent variables to update the parameters in the subsequent M-step.

Probabilistic PCA

M-step (Maximization Step):

In the M-step, we update the parameters of the model to maximize the expected log-likelihood of the observed data under the current posterior distribution of the latent variables.

The parameters to be updated include:

- The principal components matrix \mathbf{W}
- The mean vector $\boldsymbol{\mu}$
- The noise variance σ^2

The updates can be obtained by maximizing the expected log-likelihood function, which is a function of the observed data and the expected values of the latent variables.

This iterative process continues until convergence, where the changes in the parameters between successive iterations are below a predefined threshold, or when the log-likelihood of the observed data reaches a plateau.

Probabilistic PCA

6. Projection:

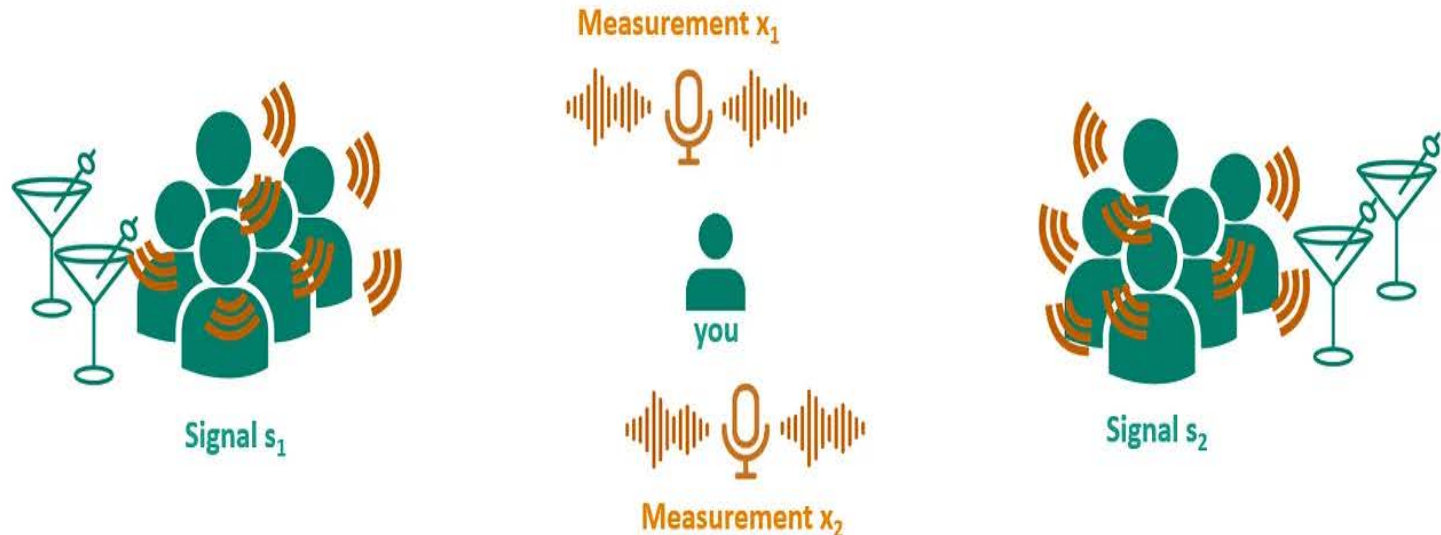
Once the parameters are estimated, the data can be projected onto the subspace spanned by the principal components \mathbf{W} . This allows for dimensionality reduction while preserving most of the variance in the data.

7. Model Evaluation:

The performance of the PPCA model can be evaluated using various criteria, such as the log-likelihood of the observed data under the model, the reconstruction error, or the ability of the model to capture the underlying structure of the data.

Independent Component Analysis(ICA)

- Independent Component Analysis is one of the various unsupervised learning algorithms which means that we do not need to supervise the model before we can use it. The origin of this method comes from signal processing where we try to separate a multivariate signal into additive subcomponents.
- The objective of ICA is to recover the original, unknown signals by separating the mixed data. The ultimate aim is to reconstruct the data such that each dimension is mutually independent.
- To make this concept more tangible, the most well-known example of ICA, the “cocktail party problem,” will be utilized.



Independent Component Analysis(ICA)

First measurement:

$$x_1 = a_{11} * s_1 + a_{12} * s_2$$

Second measurement:

$$x_2 = a_{21} * s_1 + a_{22} * s_2$$

General framework:

$$\vec{x} = A * \vec{s}$$

where a_{ij} are mixing coefficients

Goal:

$$\vec{s} = A^{-1} * \vec{x}$$

independent signals

Independent Component Analysis(ICA)

General framework:

$$\vec{x} = A * \vec{s}$$

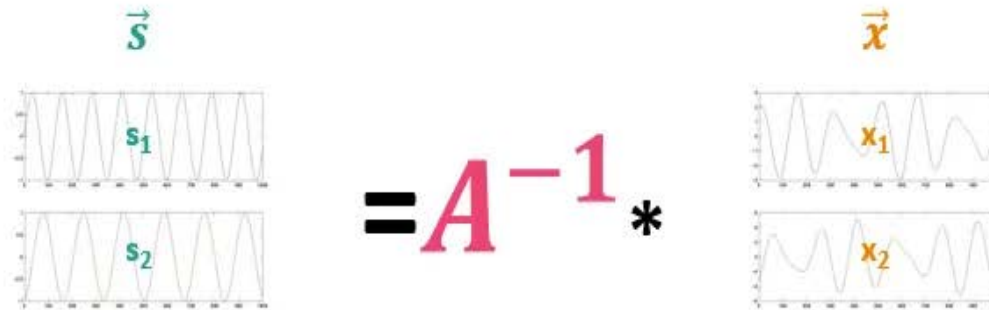
where a_{ij} are mixing coefficients



Goal:

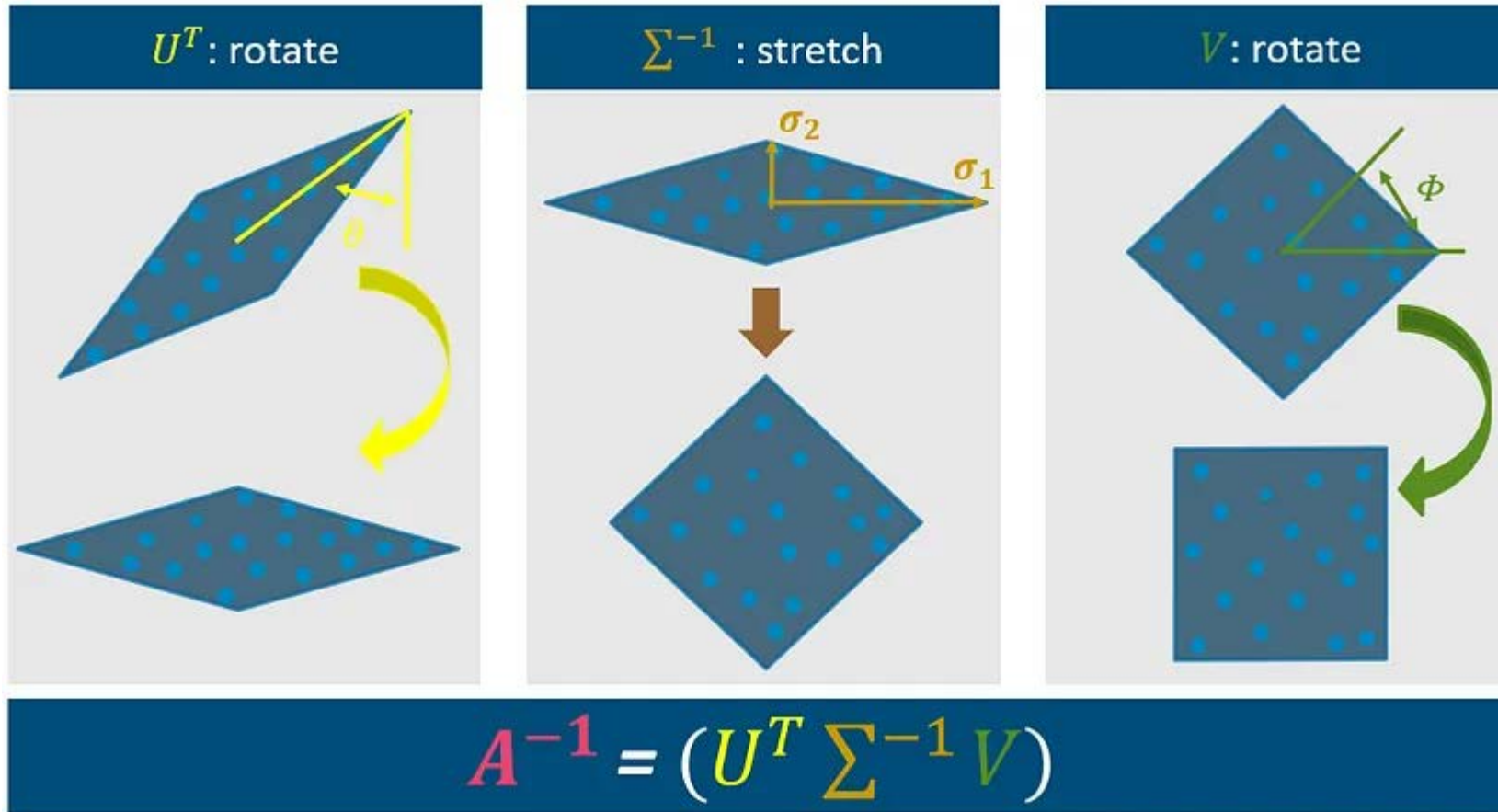
$$\vec{s} = A^{-1} * \vec{x}$$

independent signals



Independent Component Analysis(ICA)

- Separation process | the 3-step-ICA-algorithm



Independent Component Analysis(ICA)

- **Step 1:** Find the angle with maximal variance to rotate | estimate \mathbf{U}^T

The first component of the algorithm involves the use of the matrix \mathbf{U}^T , which is based on the first angle θ . The angle Theta can be derived from the primary direction of the data, as determined through Principal Component Analysis (PCA). This step rotates the graph to the position shown above.

- **Step 2:** Find the scaling of the principal components | estimate $\Sigma^{(-1)}$

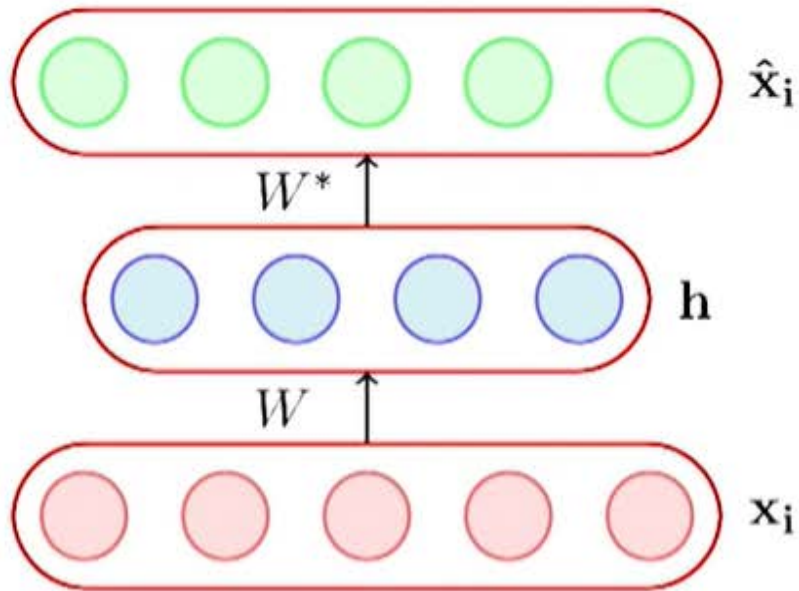
The second component involves stretching the figure, which is achieved through the Σ^{-1} step. This step employs the variances of sigma 1 and sigma 2 from the data, similar to the approach utilized in PCA.

- **Step 3:** Independence and kurtosis assumptions for rotation | estimate \mathbf{V}

The final component, which distinguishes the current algorithm from PCA, involves the rotation of the signals around angle Phi. This step aims to rebuild the original dimensions of the signals by utilizing the independence and kurtosis assumptions for rotation.

In summary, the algorithm employs measurements and performs rotation around theta, stretching through the use of variances sigma 1 and 2, and finally, rotation around Phi.

Auto Encoder



$$\mathbf{h} = g(W\mathbf{x}_i + \mathbf{b})$$

$$\hat{\mathbf{x}}_i = f(W^*\mathbf{h} + \mathbf{c})$$

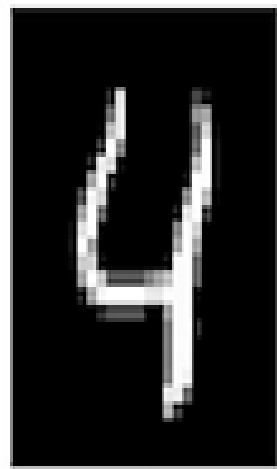
- An autoencoder is a special type of feed forward neural network which does the following
- Encodes its input \mathbf{x}_i into a hidden representation \mathbf{h}
- Decodes the input again from this hidden representation
- The model is trained to minimize a certain loss function which will ensure that $\hat{\mathbf{x}}_i$ is close to \mathbf{x}_i (we will see some such loss functions soon)

Auto Encoder

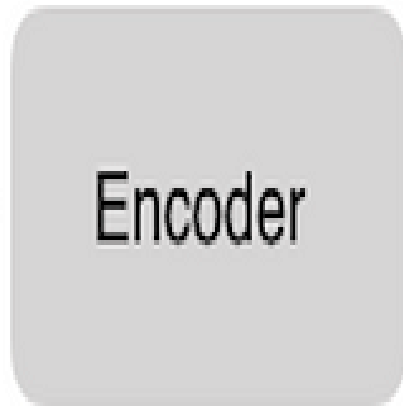
Original Input

Latent Representation

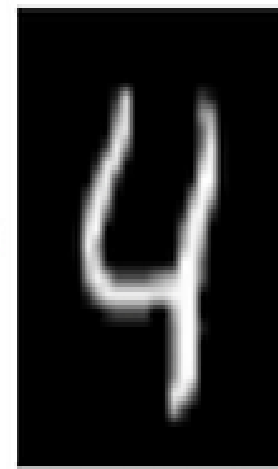
Reconstructed Output



x

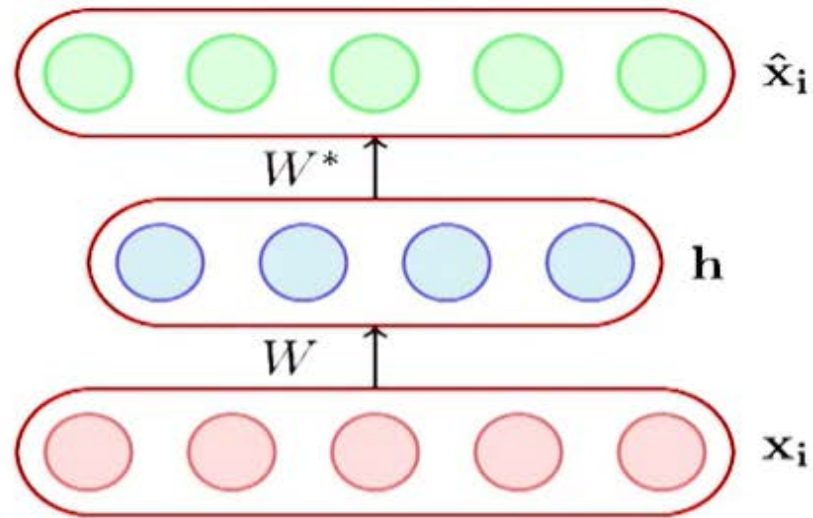


h



r

Auto Encoder

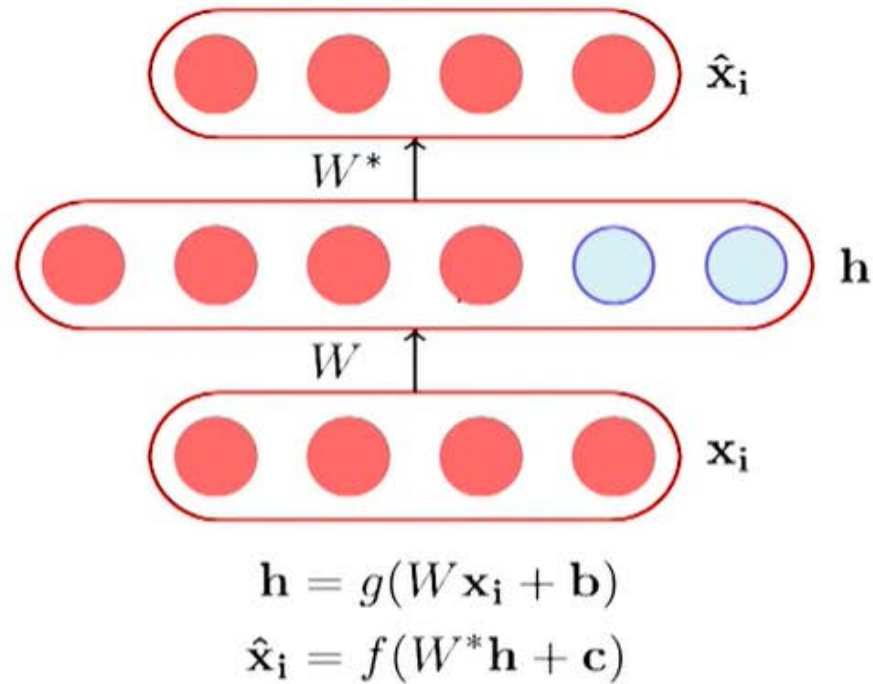


$$\mathbf{h} = g(W\mathbf{x}_i + \mathbf{b})$$
$$\hat{\mathbf{x}}_i = f(W^*\mathbf{h} + \mathbf{c})$$

- Let us consider the case where $\dim(\mathbf{h}) < \dim(\mathbf{x}_i)$
- If we are still able to reconstruct $\hat{\mathbf{x}}_i$ perfectly from \mathbf{h} , then what does it say about \mathbf{h} ?
- \mathbf{h} is a loss-free encoding of \mathbf{x}_i . It captures all the important characteristics of \mathbf{x}_i
- Do you see an analogy with PCA?

An autoencoder where $\dim(\mathbf{h}) < \dim(\mathbf{x}_i)$ is called an under complete autoencoder

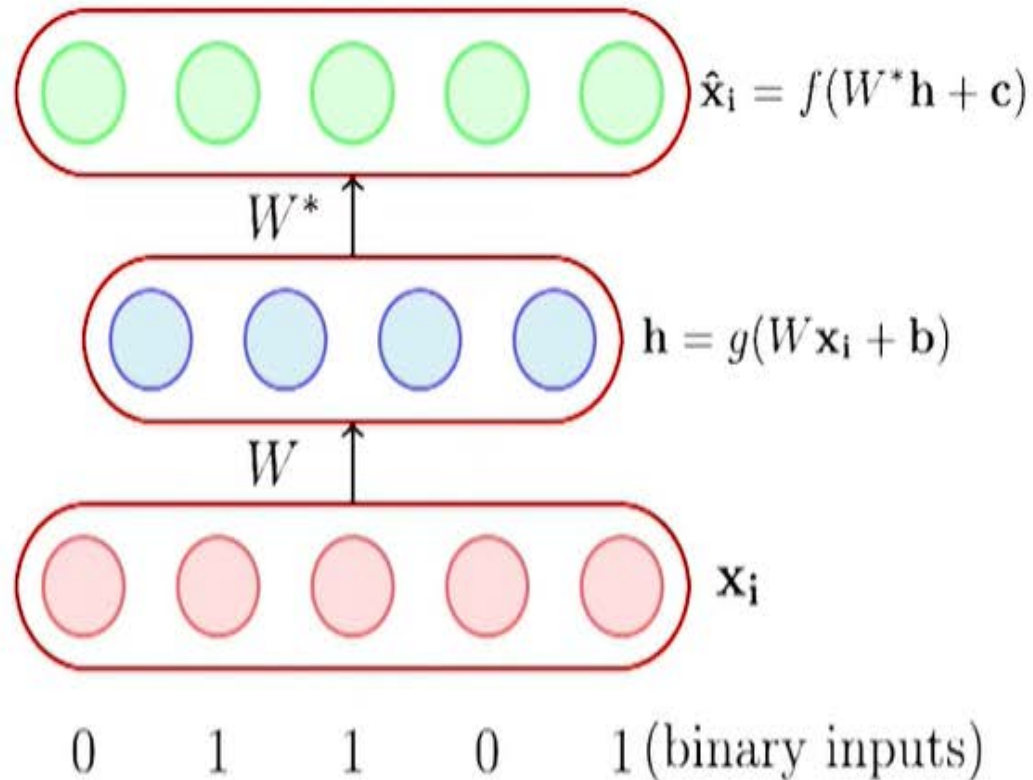
Auto Encoder



- Let us consider the case when $\dim(\mathbf{h}) \geq \dim(\mathbf{x}_i)$
- In such a case the autoencoder could learn a trivial encoding by simply copying \mathbf{x}_i into \mathbf{h} and then copying \mathbf{h} into $\hat{\mathbf{x}}_i$
- Such an identity encoding is useless in practice as it does not really tell us anything about the important characteristics of the data

An autoencoder where $\dim(\mathbf{h}) \geq \dim(\mathbf{x}_i)$ is called an over complete autoencoder

Auto Encoder



- Suppose all our inputs are binary (each $x_{ij} \in \{0, 1\}$)
- Which of the following functions would be most apt for the decoder?

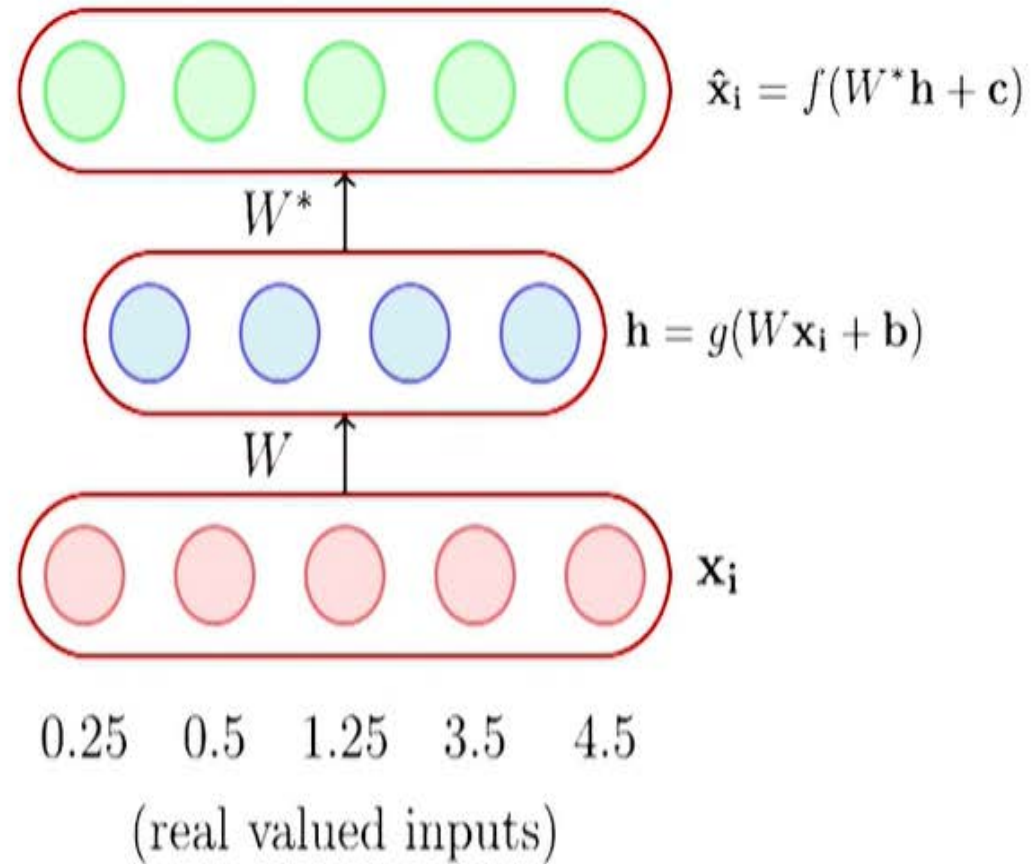
$$\hat{\mathbf{x}}_i = \tanh(W^*\mathbf{h} + \mathbf{c})$$

$$\hat{\mathbf{x}}_i = W^*\mathbf{h} + \mathbf{c}$$

$$\hat{\mathbf{x}}_i = \text{logistic}(W^*\mathbf{h} + \mathbf{c})$$

- Logistic as it naturally restricts all outputs to be between 0 and 1

Auto Encoder



- Suppose all our inputs are real (each $x_{ij} \in \mathbb{R}$)
- Which of the following functions would be most apt for the decoder?

$$\hat{\mathbf{x}}_i = \tanh(W^*\mathbf{h} + \mathbf{c})$$

$$\hat{\mathbf{x}}_i = W^*\mathbf{h} + \mathbf{c}$$

$$\hat{\mathbf{x}}_i = \text{logistic}(W^*\mathbf{h} + \mathbf{c})$$

Auto encoder- Sparse Auto encoder

The architecture of sparse autoencoders is similar to standard autoencoders, consisting of an encoder and a decoder. However, sparse autoencoders incorporate a sparsity constraint to encourage the activations of the hidden neurons to be sparse. Here's a high-level description of the architecture along with equations:

1. Encoder:

- The encoder maps the input data x to the hidden representation h using a set of weights W and biases b . The activation function f is typically a nonlinear function such as sigmoid or ReLU.
- The hidden representation h is then passed through an activation function f to introduce nonlinearity.
- The output of the encoder is the sparse representation h .

Mathematically:

$$h = f(Wx + b)$$

Auto encoder- Sparse Autoencoder

2. Decoder:

- The decoder reconstructs the input data x' from the hidden representation h using another set of weights W' and biases b' .
- The reconstruction is achieved by applying an activation function g to the linear transformation $W'h + b'$.

Mathematically:

$$x' = g(W'h + b')$$

3. Sparsity Constraint:

- In addition to minimizing the reconstruction error between the input x and the output x' , sparse autoencoders introduce a sparsity constraint on the activations of the hidden neurons.
- This is typically achieved by adding a regularization term to the loss function that penalizes non-zero activations. One common choice is the L1 regularization term.

Mathematically, the sparsity term can be defined as:

$$\Omega(h) = \lambda \sum_j |h_j|$$

where λ is the sparsity weight and $|h_j|$ is the absolute value of the activation of the j th neuron in the hidden layer.

Auto encoder- Sparse Autoencoder

The overall loss function for sparse autoencoders combines the reconstruction error term and the sparsity term:

$$L(x, x') = \|x - x'\|^2 + \Omega(h)$$

The model is then trained to minimize this combined loss function using optimization techniques such as gradient descent.

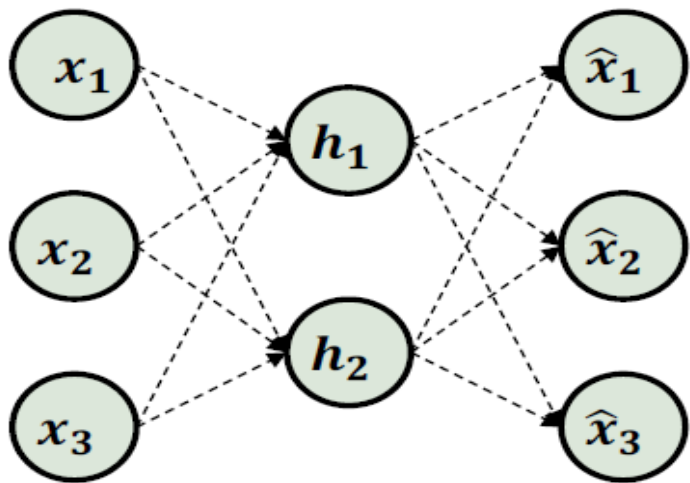
In these equations:

- x represents the input data.
- h represents the hidden representation (sparse encoding).
- x' represents the reconstructed output.
- W and W' are the weight matrices connecting the input to the hidden layer and the hidden layer to the output layer, respectively.
- b and b' are the bias vectors.
- f and g are activation functions.
- λ is the sparsity weight, controlling the impact of the sparsity constraint on the overall loss function.



Sparse autoencoders

- If the capacity of the autoencoder is too large, it can learn identity function
- Add L1 regularization term to constrain capacity

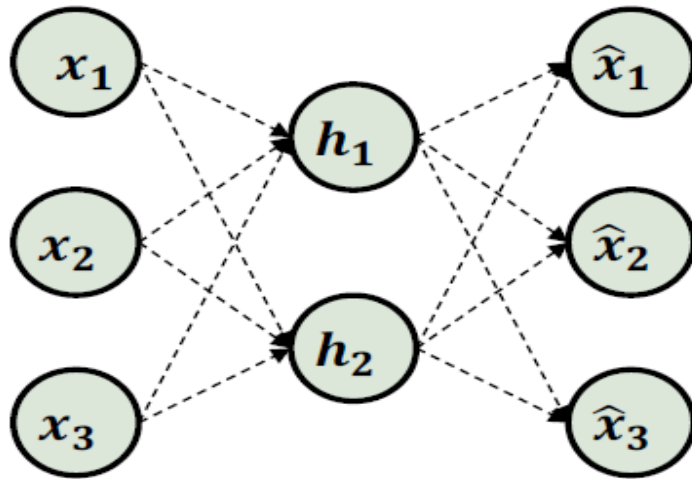


$$L(\mathbf{x}, g(f(\mathbf{x}))) + \Omega(\mathbf{h})$$

$$\Omega(\mathbf{h}) = \lambda \sum_i |h_i|$$

Contractive autoencoder

- Sparse autoencoder that replaces penalty term to differentiation of h
- This make autoencoder robust to small perturbation of input data x



$$L(x, g(f(x))) + \Omega(h, x),$$

$$\Omega(h, x) = \lambda \sum_i \|\nabla_x h_i\|^2.$$

Regularizing by Penalizing Derivatives

- Another strategy for regularizing an autoencoder
- Use penalty as in sparse autoencoders

$$L(\mathbf{x}, g(f(\mathbf{x}))) + \Omega(\mathbf{h}, \mathbf{x})$$

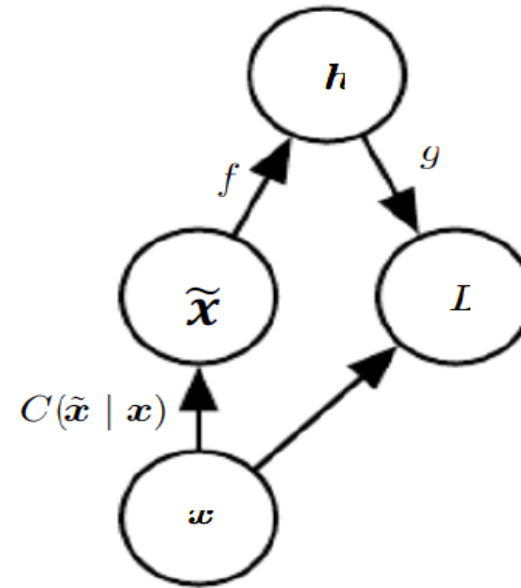
- But with a different form of Ω

$$\Omega(\mathbf{h}, \mathbf{x}) = \lambda \sum_i \|\nabla_{\mathbf{x}} h_i\|^2$$

- Forces the model to learn a function that does not change much when \mathbf{x} changes slightly
- Called a *Contractive Auto Encoder* (CAE)
- This model has theoretical connections to
 - Denoising autoencoders
 - Manifold learning
 - Probabilistic modeling

Denoising autoencoders

- To make network robust about noise through blow procedure
 1. sample training example x from the training data
 2. sample a corrupted version \tilde{x} from $C(\tilde{x}|x)$
 3. use \tilde{x} as input data and x as label

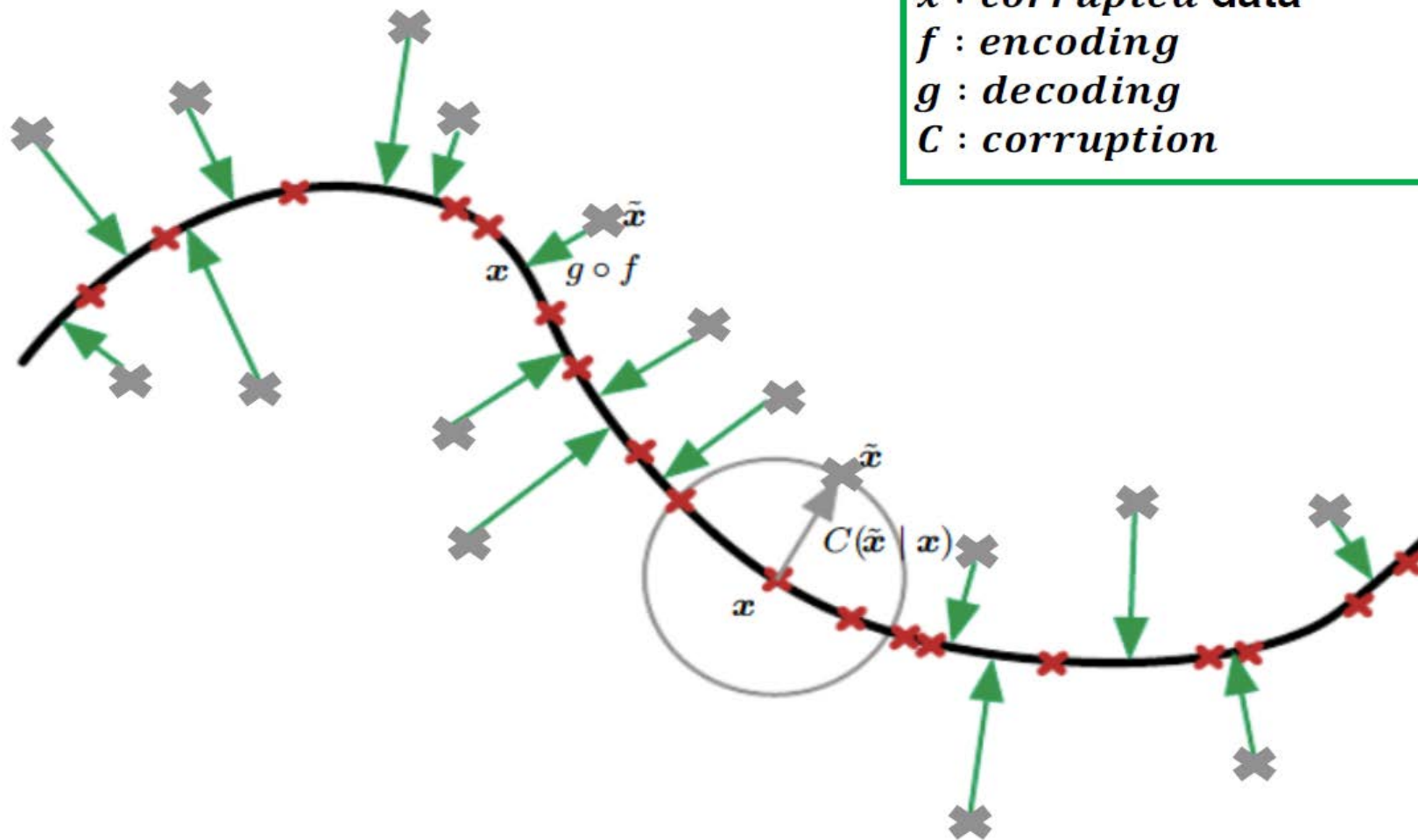


Denoising Autoencoders (DAE)

- Rather than adding a penalty Ω to the cost function, we can obtain an autoencoder that learns something useful
 - By changing the reconstruction error term of the cost function
- Traditional autoencoders minimize $L(\mathbf{x}, g(f(\mathbf{x})))$
 - where L is a loss function penalizing $g(f(\mathbf{x}))$ for being dissimilar from \mathbf{x} , such as L^2 norm of difference: mean squared error
- A DAE minimizes $L(\mathbf{x}, g(f(\tilde{\mathbf{x}})))$
 - where $\tilde{\mathbf{x}}$ is a copy of \mathbf{x} that has been corrupted by some form of noise
 - The autoencoder must undo this corruption rather than simply copying their input
- Denoising training forces f and g to implicitly learn the structure of $p_{\text{data}}(\mathbf{x})$
- Another example of how useful properties can emerge as a by-product of minimizing reconstruction error

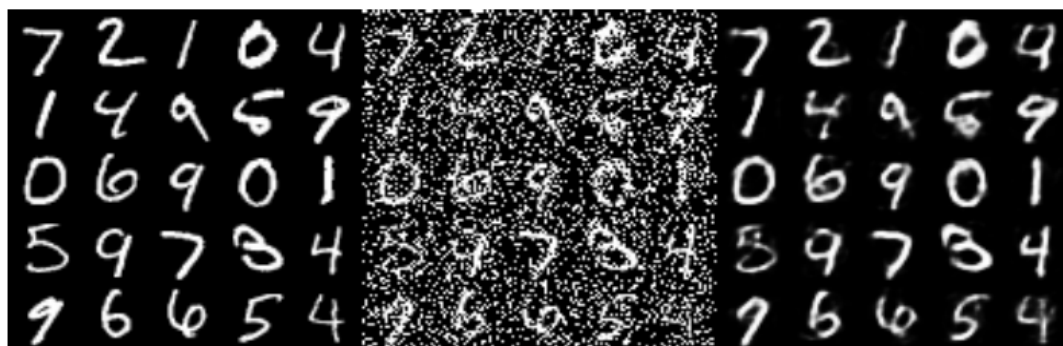
Denoising autoencoders

x : original data
 \tilde{x} : corrupted data
 f : encoding
 g : decoding
 C : corruption



Denoising autoencoders

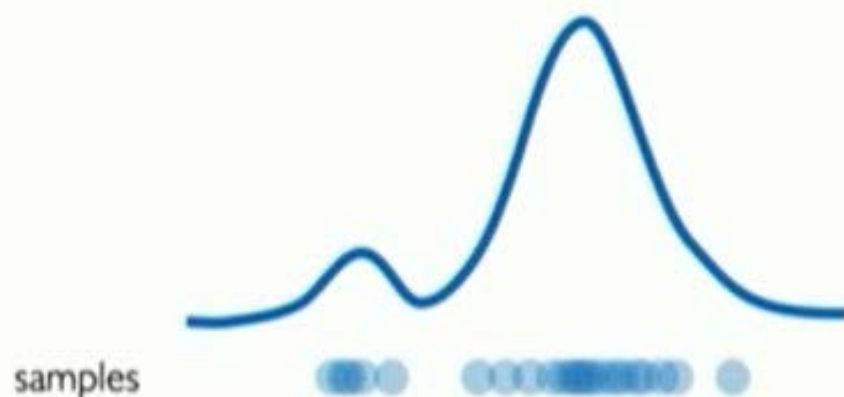
- An example about denoise through denoising autoencoder



Generative modeling

Goal: Take as input training samples from some distribution and learn a model that represents that distribution

Density Estimation



Sample Generation



Input samples

Training data $\sim P_{data}(x)$



Generated samples

Generated $\sim P_{model}(x)$

How can we learn $P_{model}(x)$ similar to $P_{data}(x)$?

Why generative models? Debiasing

Capable of uncovering **underlying features** in a dataset



Homogeneous skin color, pose

VS



Diverse skin color, pose, illumination

How can we use this information to create fair and representative datasets?

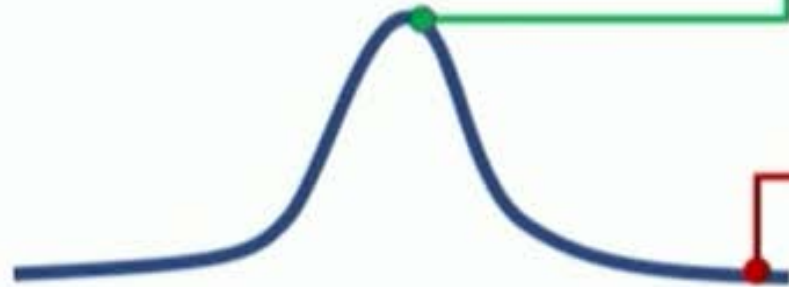
Why generative models? Outlier detection

- **Problem:** How can we detect when we encounter something new or rare?
- **Strategy:** Leverage generative models, detect outliers in the distribution
- Use outliers during training to improve even more!

95% of Driving Data:
(1) sunny, (2) highway, (3) straight road



Detect outliers to avoid unpredictable behavior when training



Edge Cases



Harsh Weather



Pedestrians

Latent variable models

Autoencoders and Variational
Autoencoders (VAEs)



Generative Adversarial
Networks (GANs)



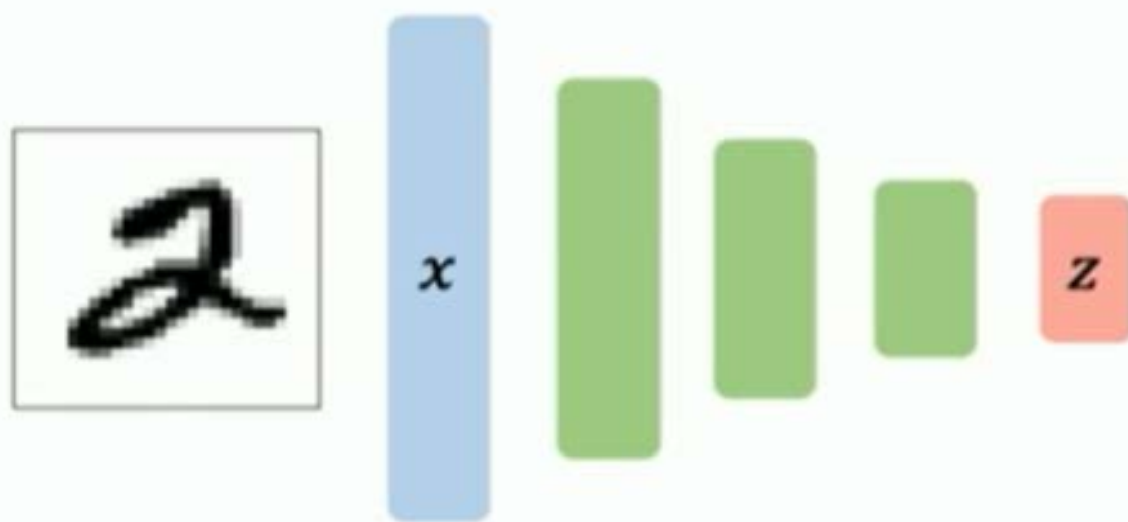
What is a latent variable?



Myth of the Cave

Autoencoders: background

Unsupervised approach for learning a **lower-dimensional** feature representation from unlabeled training data

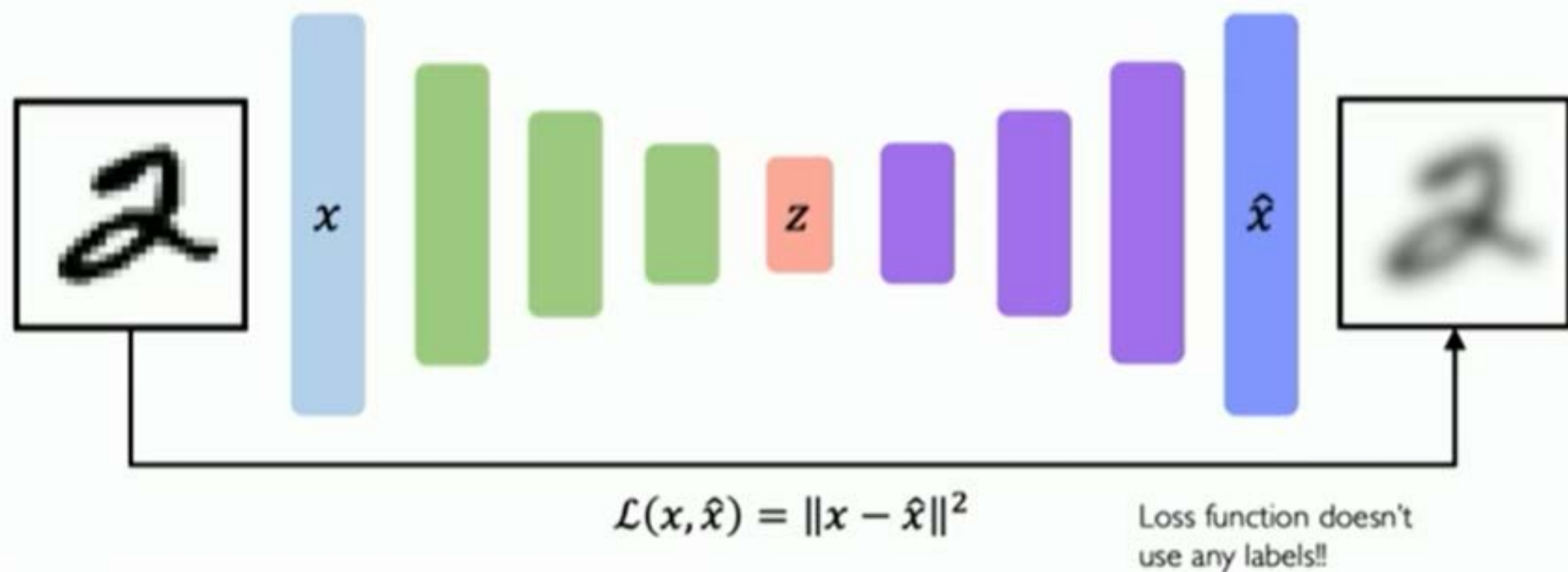


"Encoder" learns mapping from the data, x , to a low-dimensional latent space, z

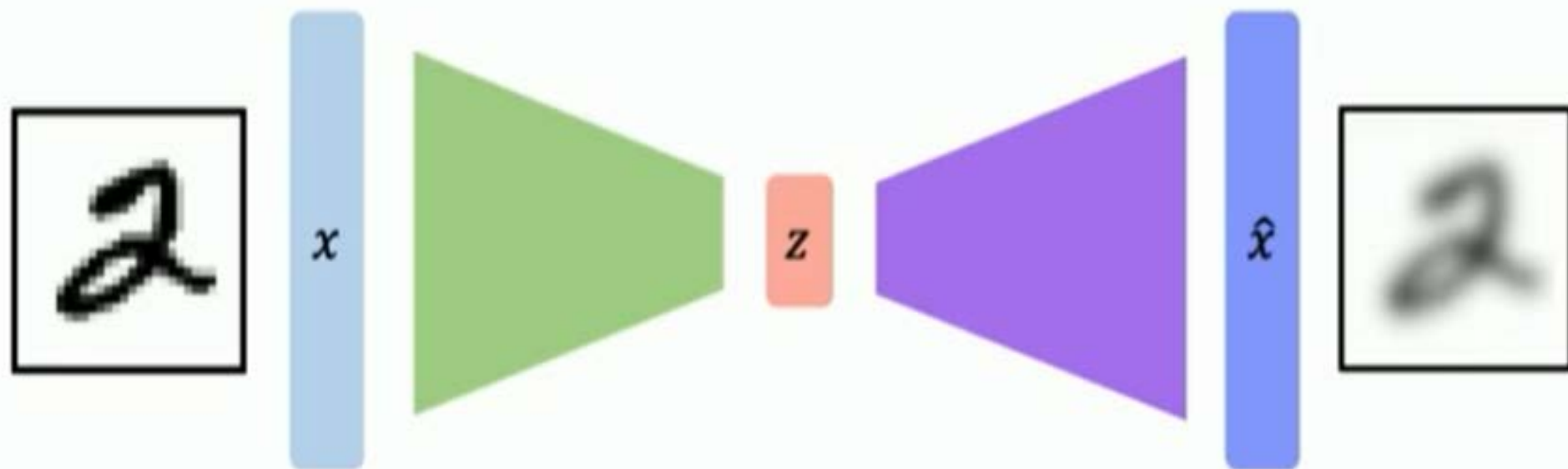
Autoencoders: background

How can we learn this latent space?

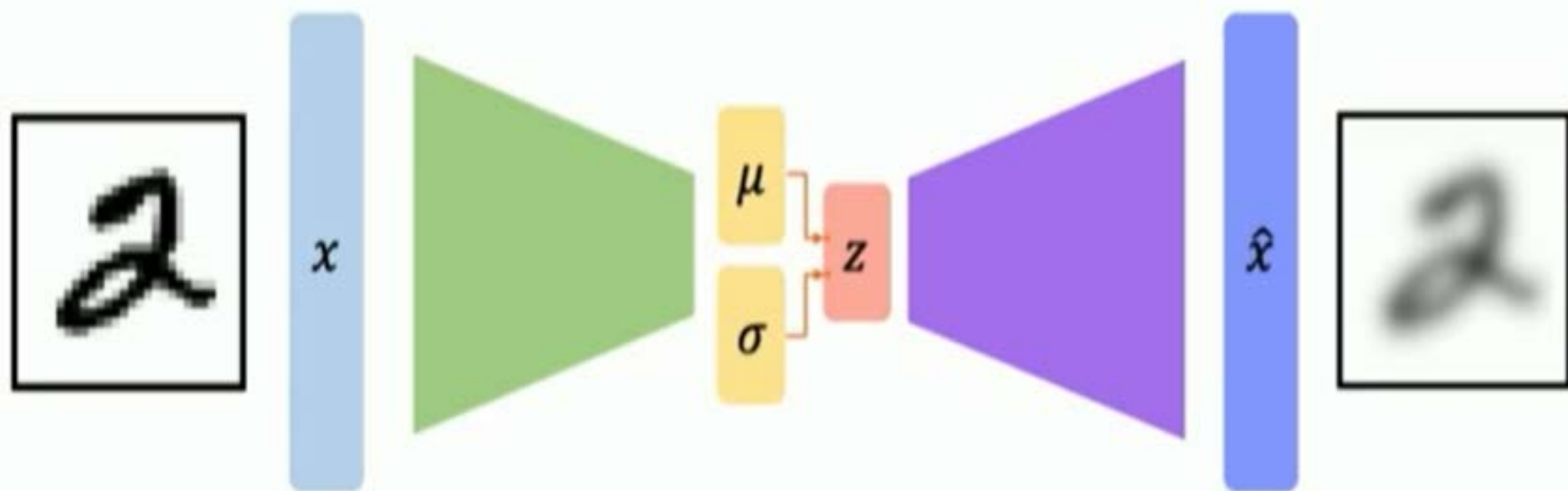
Train the model to use these features to **reconstruct the original data**



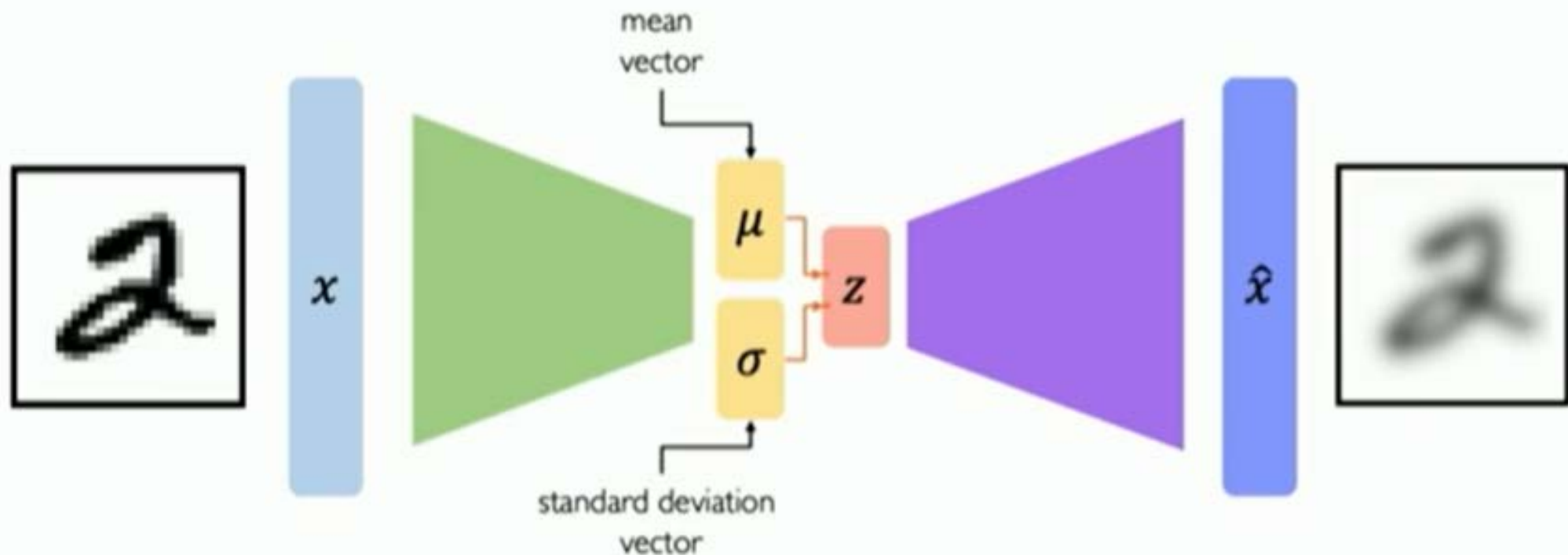
Traditional autoencoders



VAEs: key difference with traditional autoencoder



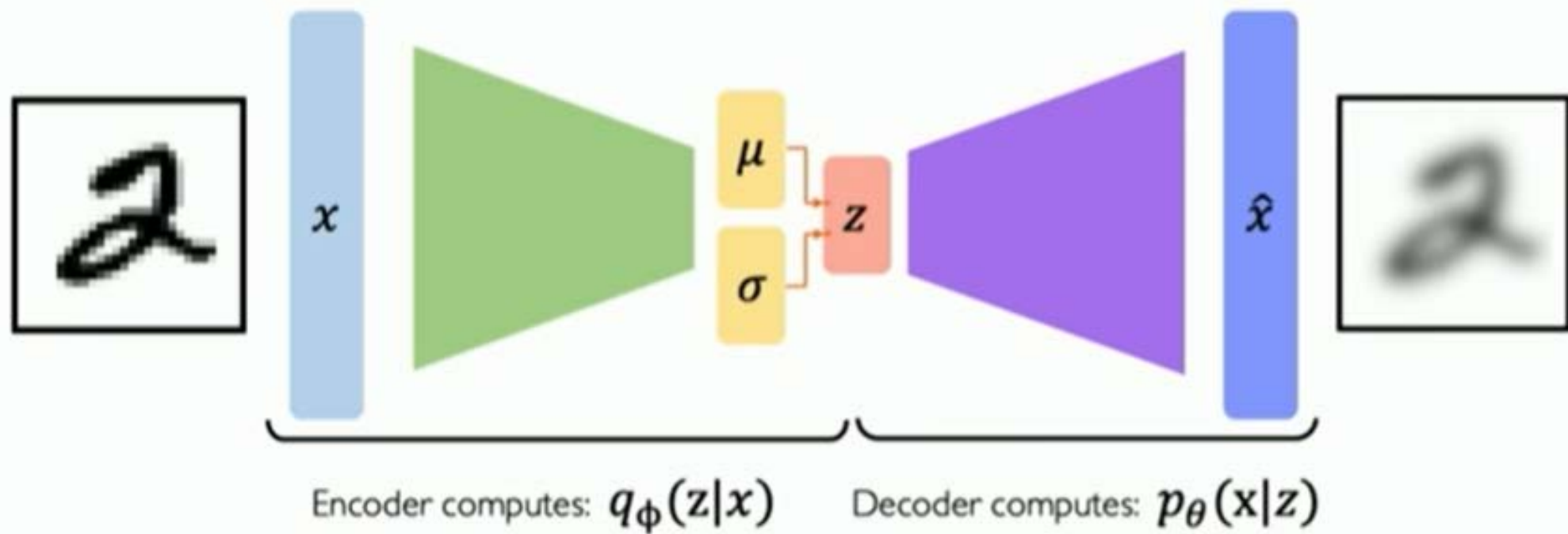
VAEs: key difference with traditional autoencoder



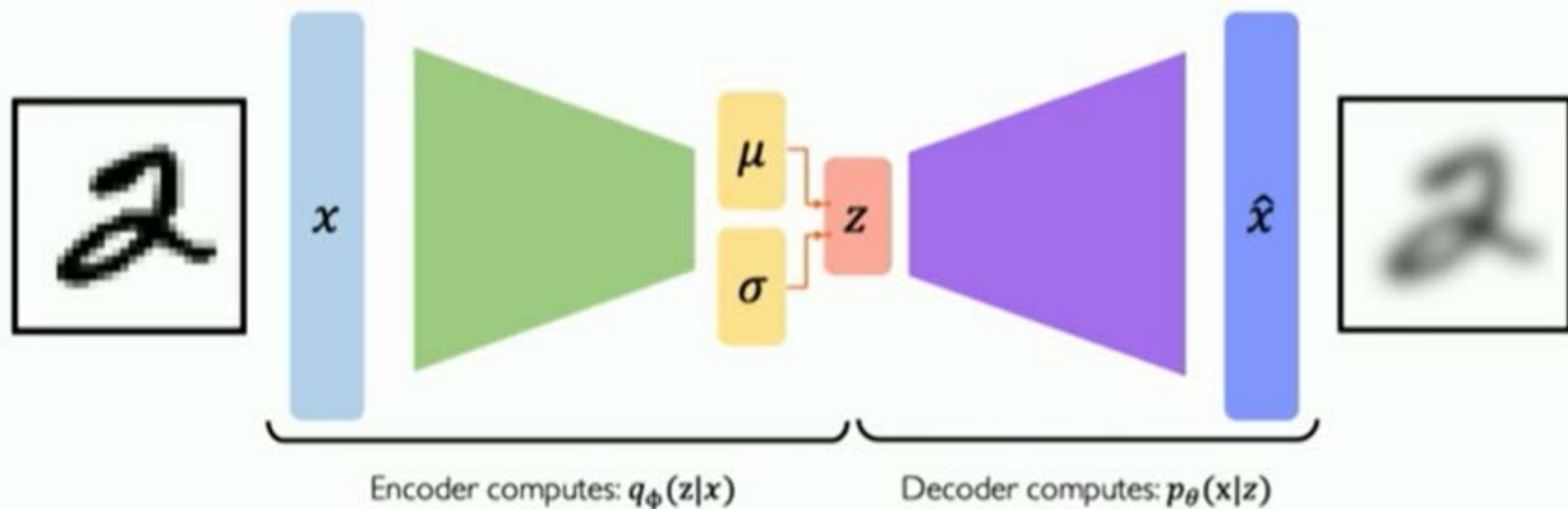
Variational autoencoders are a probabilistic twist on autoencoders!

Sample from the mean and standard deviation to compute latent sample

VAE optimization

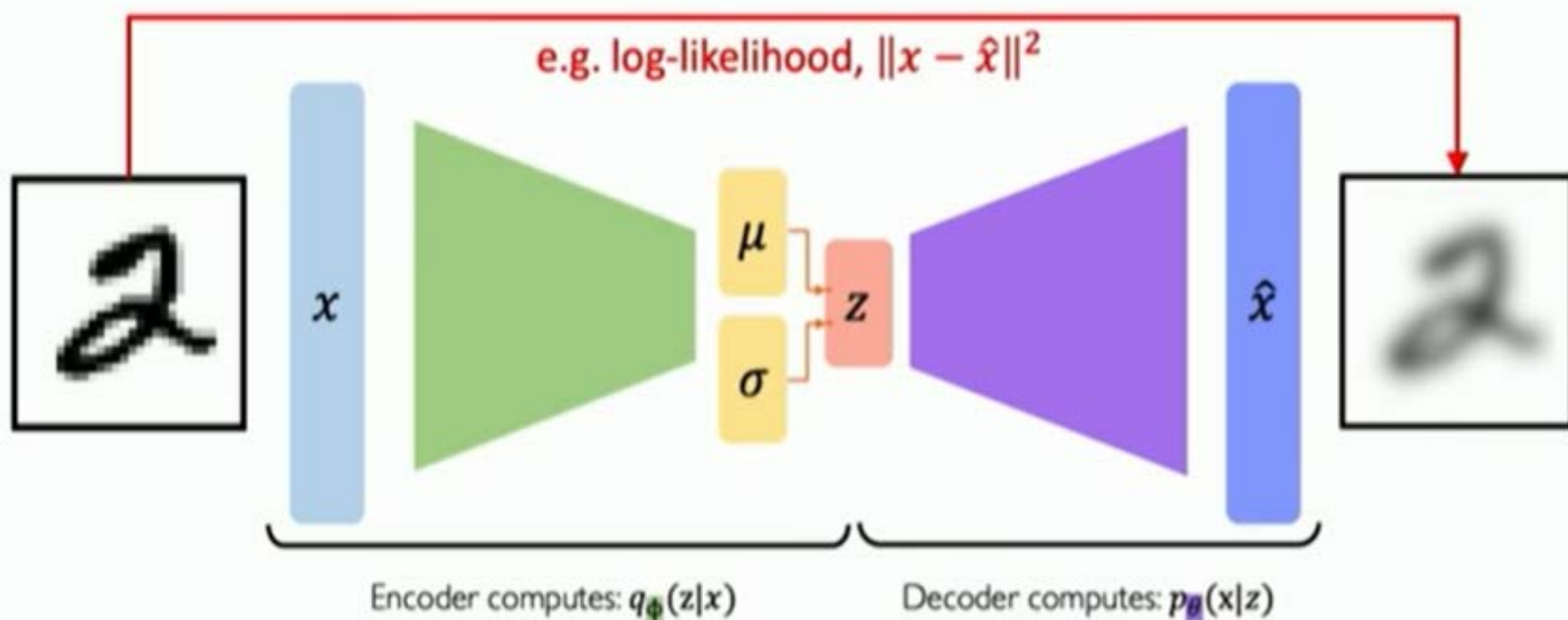


VAE optimization



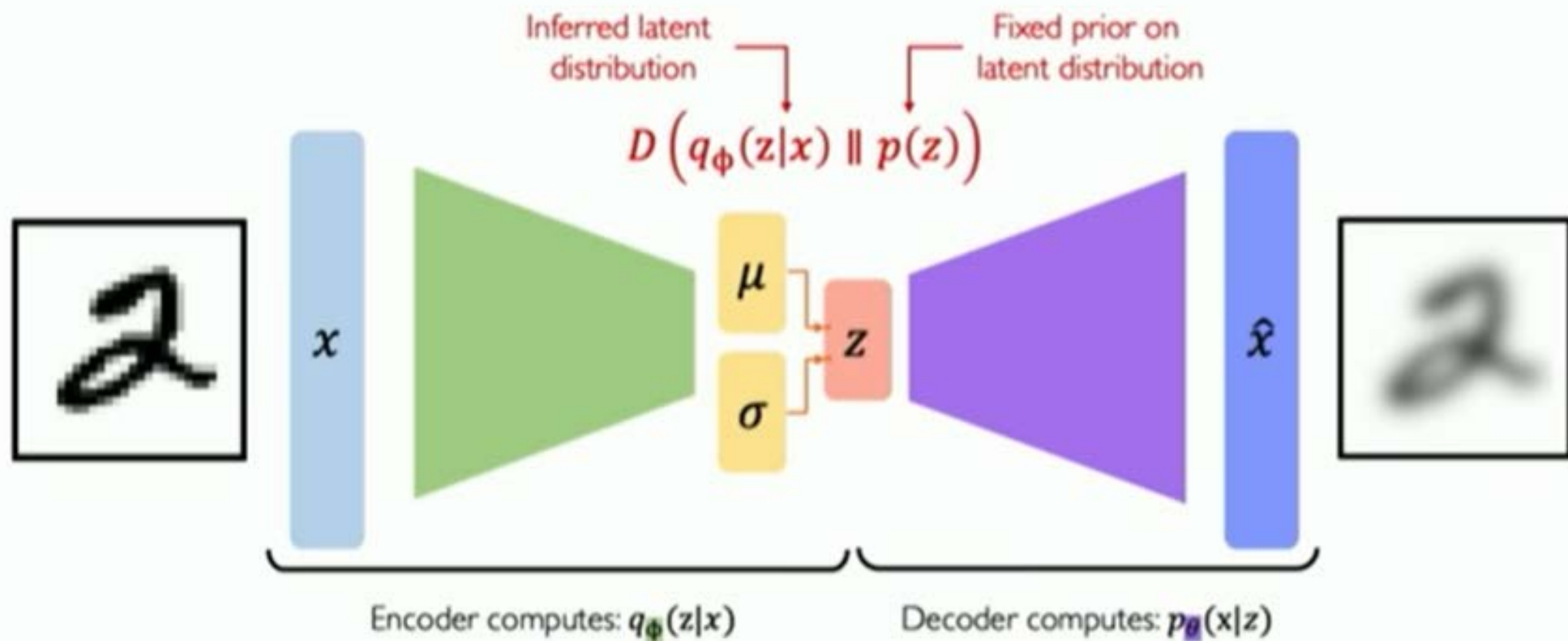
$$\mathcal{L}(\phi, \theta, x) = (\text{reconstruction loss}) + (\text{regularization term})$$

VAE optimization



$$\mathcal{L}(\phi, \theta, x) = (\text{reconstruction loss}) + (\text{regularization term})$$

VAE optimization



$$\mathcal{L}(\phi, \theta, x) = (\text{reconstruction loss}) + (\text{regularization term})$$

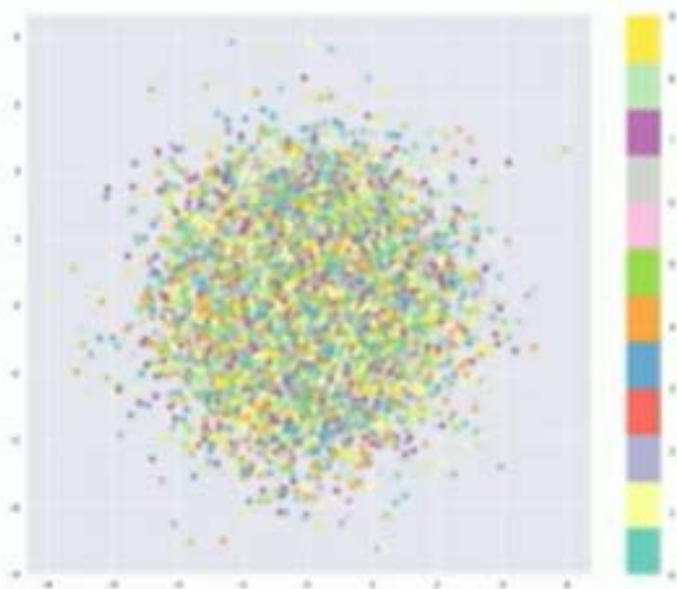
Priors on the latent distribution

$$D \left(q_{\phi}(z|x) \parallel p(z) \right)$$

Inferred latent
distribution



Fixed prior on
latent distribution



Common choice of prior – Normal Gaussian:

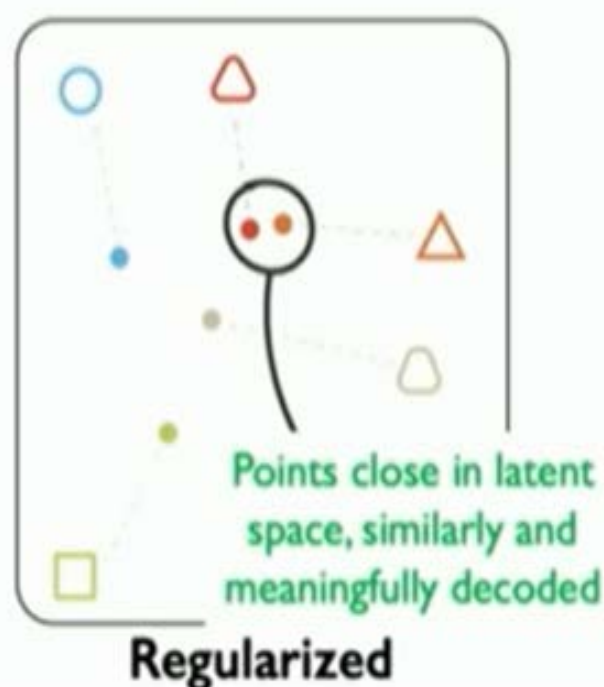
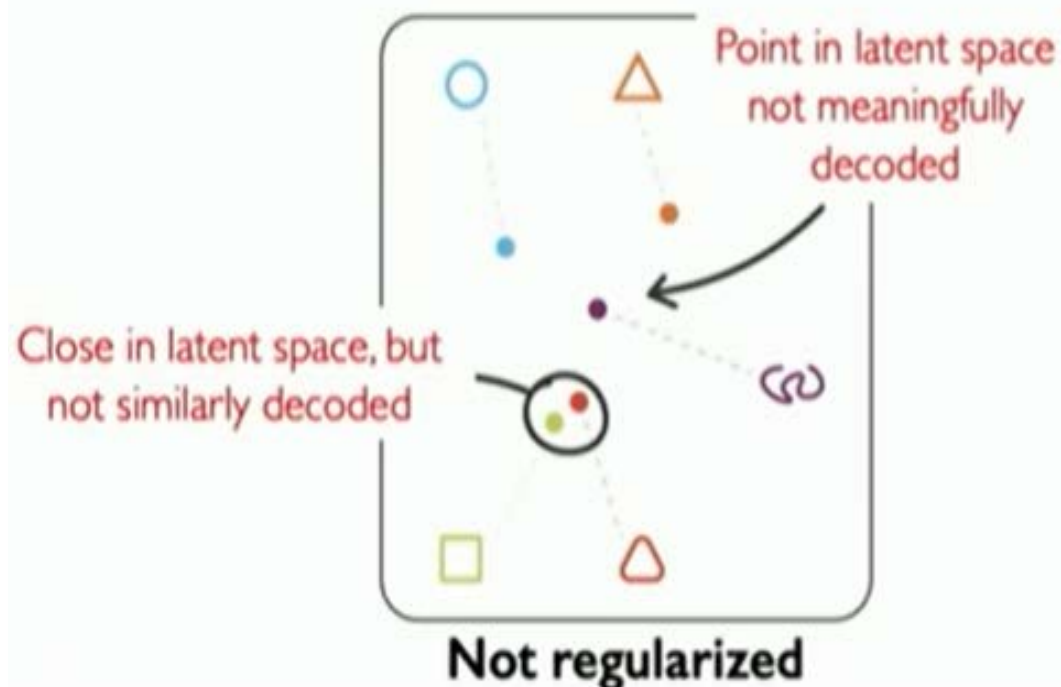
$$p(z) = \mathcal{N}(\mu = 0, \sigma^2 = 1)$$

- Encourages encodings to distribute evenly around the center of the latent space
- Penalize the network when it tries to "cheat" by clustering points in specific regions (i.e., by memorizing the data)

Intuition on regularization and the Normal prior

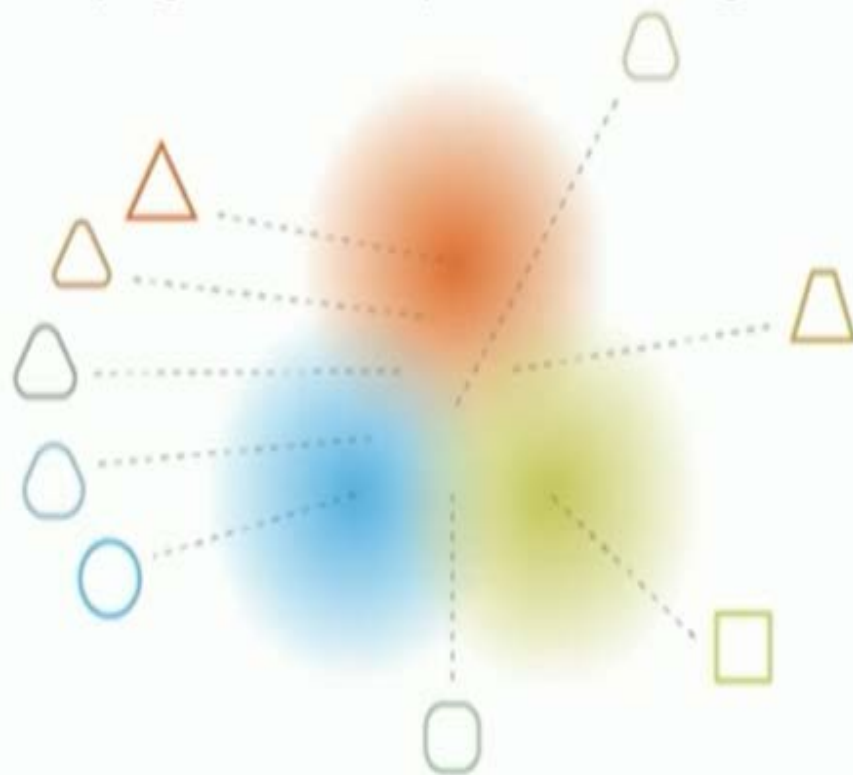
What properties do we want to achieve from regularization? 🤔

1. **Continuity:** points that are close in latent space \rightarrow similar content after decoding
2. **Completeness:** sampling from latent space \rightarrow "meaningful" content after decoding



Intuition on regularization and the Normal prior

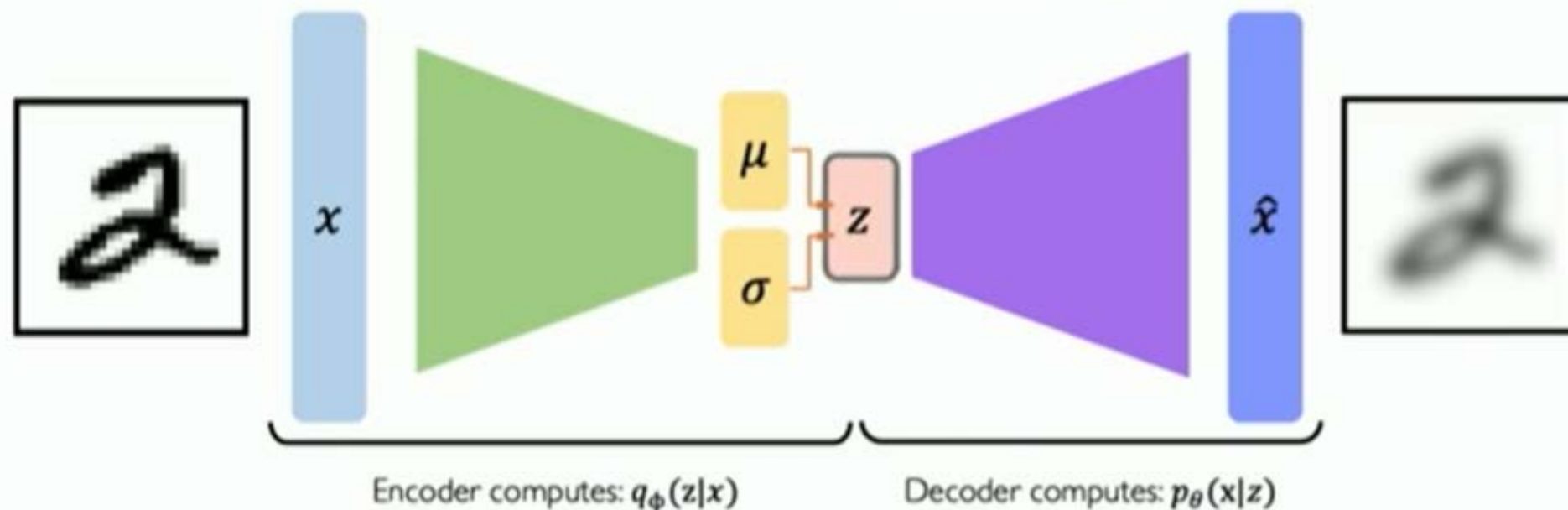
1. **Continuity**: points that are close in latent space \rightarrow similar content after decoding
2. **Completeness**: sampling from latent space \rightarrow "meaningful" content after decoding



Regularization with Normal prior helps enforce **information gradient** in the latent space.

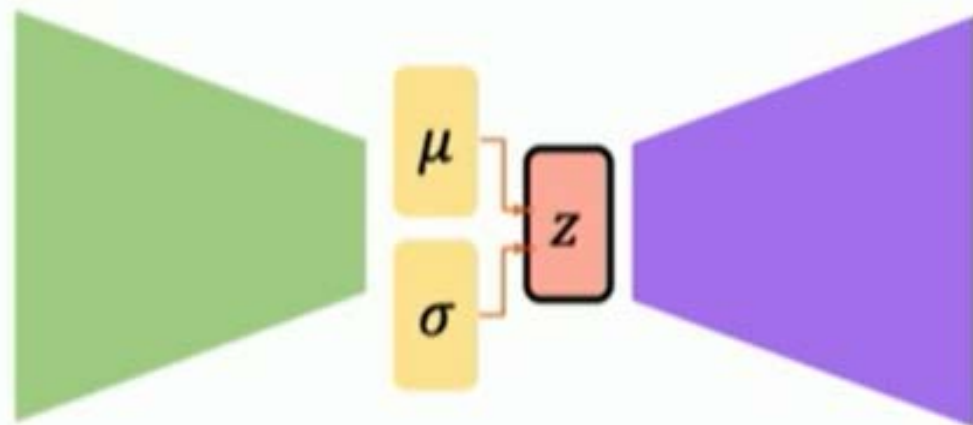
VAE computation graph

Problem: We cannot backpropagate gradients through sampling layers!



$$\mathcal{L}(\phi, \theta, x) = (\text{reconstruction loss}) + (\text{regularization term})$$

Reparametrizing the sampling layer



Key Idea:

~~$z \sim \mathcal{N}(\mu, \sigma^2)$~~

Consider the sampled latent vector z as a sum of

- a fixed μ vector,
- and fixed σ vector, scaled by random constants drawn from the prior distribution

$$\Rightarrow z = \mu + \sigma \odot \epsilon$$

where $\epsilon \sim \mathcal{N}(0,1)$

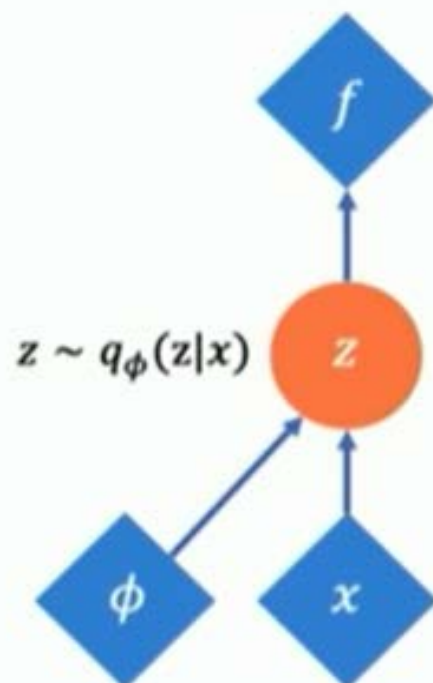
Reparametrizing the sampling layer



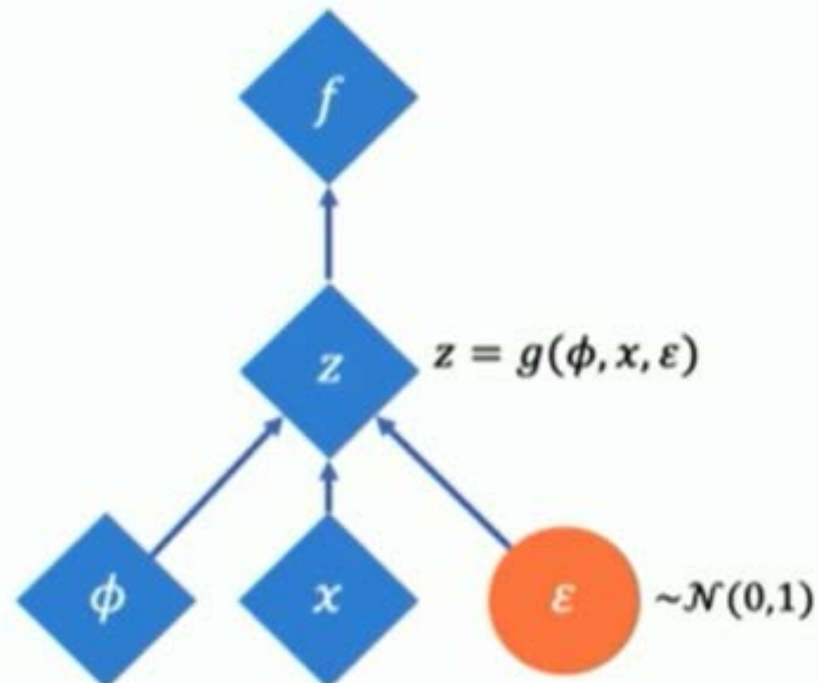
Deterministic node



Stochastic node



Original form



Reparametrized form

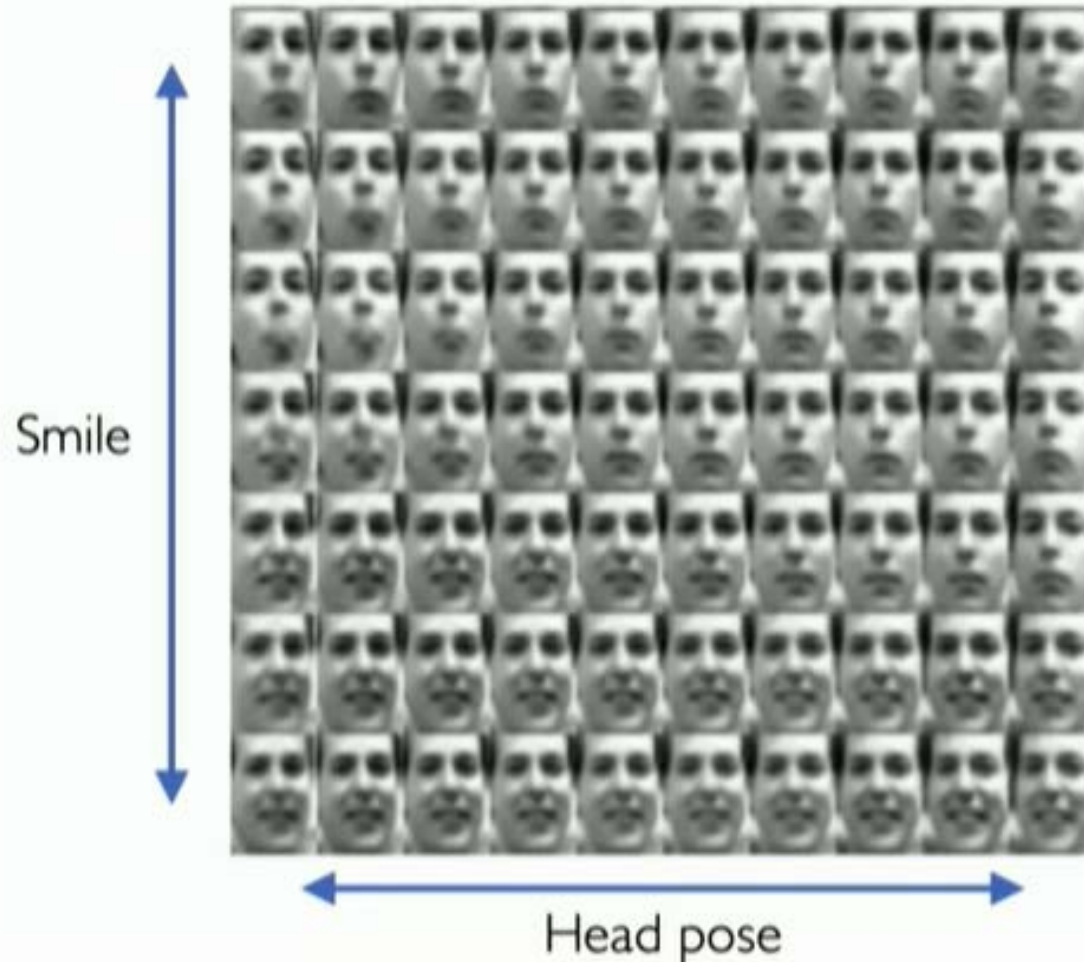
VAEs: Latent perturbation

Slowly increase or decrease a **single latent variable**
Keep all other variables fixed



Head pose

VAEs: Latent perturbation



Ideally, we want latent variables that are uncorrelated with each other

Enforce diagonal prior on the latent variables to encourage independence

Disentanglement

Latent space disentanglement with β -VAEs

β -VAE loss:

$$\mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z}, \beta) = \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction term}} - \underbrace{\beta D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}))}_{\text{Regularization term}}$$

$\beta > 1$: constrain latent bottleneck, encourage efficient latent encoding \rightarrow disentanglement

Head rotation (azimuth)



Standard VAE ($\beta = 1$)

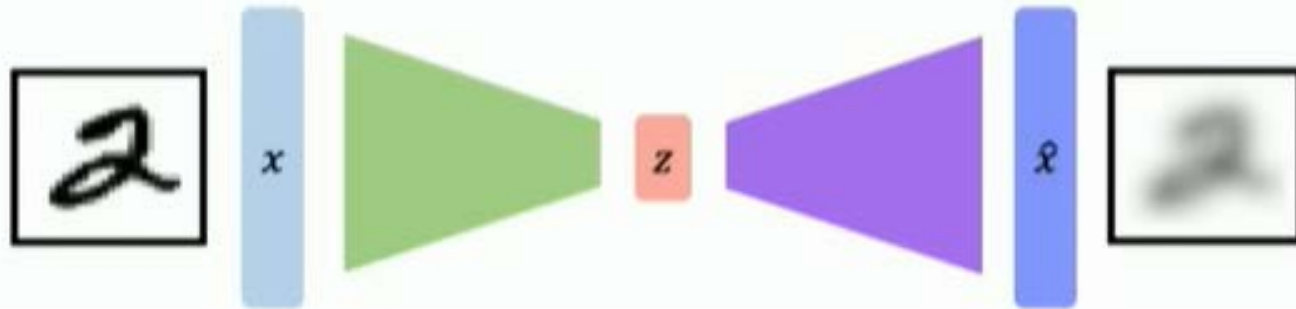


β -VAE ($\beta = 250$)

Smile also
changing!

VAE summary

1. Compress representation of world to something we can use to learn
2. Reconstruction allows for unsupervised learning (no labels!)
3. Reparameterization trick to train end-to-end
4. Interpret hidden latent variables using perturbation
5. Generating new examples



Applications of Auto encoders

Autoencoders are a type of artificial neural network used for unsupervised learning. They learn how to efficiently encode and then decode data, typically by minimizing the reconstruction error between the input and the output. Here are some detailed applications:

➤ Dimensionality Reduction:

- One of the primary applications of autoencoders is reducing the dimensionality of data. By training an autoencoder to reconstruct high-dimensional data from a lower-dimensional representation, you effectively learn a compressed representation of the data. This is useful for data visualization, feature extraction, and speeding up subsequent processing steps.

➤ Anomaly Detection:

- Autoencoders can be used for anomaly detection in data. By training the autoencoder on normal data, it learns to reconstruct typical instances accurately. When presented with anomalous data, the reconstruction error is typically higher, signaling an anomaly. This makes autoencoders effective for detecting outliers or anomalies in various domains such as cybersecurity, fraud detection, and manufacturing quality control.

➤ Image Denoising:

- Autoencoders can be trained to remove noise from images. By presenting noisy images as input and clean images as output during training, the autoencoder learns to reconstruct the clean images from the noisy inputs. This learned ability can then be applied to denoise unseen images.

Applications of Auto encoders

➤ **Feature Learning:**

- Autoencoders are proficient at learning meaningful representations of data. In the context of feature learning, autoencoders can automatically discover useful features from raw data, which can then be utilized for downstream tasks such as classification or regression. This is particularly advantageous when dealing with high-dimensional and unstructured data like images, text, or audio.

➤ **Data Generation:**

- Variational Autoencoders (VAEs), a type of autoencoder, are capable of generating new data samples that resemble the training data. By sampling from the learned latent space, VAEs can generate diverse outputs, making them useful for tasks such as generating images, music, text, or other creative applications.

➤ **Representation Learning:**

- Autoencoders can learn rich representations of data without requiring labeled examples. These learned representations can capture complex structures and patterns in the data, making them useful for downstream tasks like transfer learning, where pre-trained autoencoder models are fine-tuned on smaller labeled datasets for specific tasks.

➤ **Recommendation Systems:**

- Autoencoders can be employed in recommendation systems to learn user or item embeddings. By representing users and items in a low-dimensional latent space, autoencoders can capture user-item interactions and preferences, enabling more effective recommendations for users.

➤ **Image Compression:**

- Autoencoders can be trained to compress images into a lower-dimensional representation and then reconstruct them. This compressed representation can be transmitted or stored more efficiently, making autoencoders useful for image compression applications.

Learning with Manifolds

Learning manifolds refers to the process of discovering and understanding the underlying structure of data, particularly in high-dimensional spaces, where data points often lie on or near lower-dimensional structures known as manifolds.

In simpler terms, it's about finding the essential patterns or shapes that describe how data is distributed within a complex space.

- **Manifold Definition:**

A manifold is a continuous and smooth geometric structure that can be embedded in a higher-dimensional space. For example, a two-dimensional surface like a sphere or a twisted torus is a manifold embedded in a three-dimensional space.

- **High-Dimensional Data:**

In many real-world scenarios, data is high-dimensional, meaning each data point is described by a large number of features or dimensions. However, not all of these dimensions may be equally relevant or informative.

- **Data Distribution:**

Despite the high dimensionality, the data often resides on or near lower-dimensional structures or manifolds. These manifolds represent the essential underlying structure of the data distribution.

Learning with Manifolds

- **Learning Manifolds:**

Learning manifolds involves discovering these lower-dimensional structures from the high-dimensional data. This process aims to capture the intrinsic geometry of the data distribution and uncover the underlying patterns or relationships between data points.

- **Dimensionality Reduction:**

One way to learn manifolds is through dimensionality reduction techniques, which aim to project high-dimensional data onto lower-dimensional spaces while preserving the essential structure of the data. By reducing the dimensionality, these techniques reveal the underlying manifold on which the data resides.

- **Unsupervised Learning:**

Learning manifolds often falls under the umbrella of unsupervised learning, where the goal is to extract meaningful representations or structures from unlabeled data. Algorithms such as autoencoders, principal component analysis (PCA), and t-distributed stochastic neighbor embedding (t-SNE) are commonly used for this purpose.

- **Applications:**

Understanding and learning manifolds have numerous applications across various domains, including data visualization, clustering, anomaly detection, generative modeling, and semi-supervised learning. By uncovering the intrinsic structure of the data, these applications can benefit from more efficient representation, better generalization, and improved performance.

Representation Learning

Representation learning is a subfield of machine learning concerned with learning efficient representations of data. Representation learning revolves around the idea of transforming raw, high-dimensional data into a more compact, meaningful, and informative representation. The aim is to capture the underlying structure and essential features of the data in a way that facilitates learning and understanding by machine learning algorithms.

- **Learning Paradigms:**

Representation learning encompasses various learning paradigms:

- **Unsupervised Learning:** In unsupervised learning, algorithms learn representations from unlabeled data. They rely solely on the inherent structure of the data to discover patterns and extract meaningful features without explicit supervision.
- **Semi-Supervised Learning:** This approach combines labeled and unlabeled data to learn representations more efficiently. It leverages both the labeled examples for guidance and the vast amount of unlabeled data available for additional information.
- **Self-Supervised Learning:** A form of unsupervised learning where the learning task is automatically generated from the data itself. For example, predicting masked words in a sentence or predicting missing parts of an image are self-supervised tasks.

Representation Learning

- **Hierarchical Representations:**

- Effective representation learning often involves learning hierarchical representations. In hierarchical representations, features are organized into multiple levels of abstraction. Lower layers of a model typically capture low-level features (e.g., edges, textures), while higher layers capture more abstract and complex features (e.g., object parts, object categories).

- **Transfer Learning and Fine-Tuning:**

- Representation learning enables transfer learning, a technique where knowledge acquired from solving one task is transferred to another related task. Pre-trained models, which have learned informative representations from large datasets, can be fine-tuned on smaller datasets for specific tasks, saving time and computational resources.

- **Domain Adaptation:**

- In real-world scenarios, data distributions may vary between different domains or environments. Representation learning techniques can facilitate domain adaptation by learning domain-invariant representations. By aligning the learned representations across domains, models can generalize better to unseen domains with minimal performance degradation.

Representation Learning

- **Disentangled Representation:**

- Disentangled representation learning aims to separate the underlying factors of variation in the data into independent and interpretable components. This enables better understanding and control over the learned representations, which is crucial for tasks like generative modeling, where manipulating specific attributes of the data (e.g., pose, style, or identity in images) is desired.

- **Applications:**

- Representation learning finds applications across various domains, including computer vision, natural language processing, speech recognition, recommender systems, and reinforcement learning. It serves as the backbone for many state-of-the-art machine learning algorithms and contributes to advancements in AI research and technology.
- Representation learning is fundamental to the success of machine learning algorithms, as it enables them to extract meaningful features and patterns from raw data, leading to improved performance on various tasks and applications.

