

Exam 3 Review Day

Contents

- expectation maximization, k-means
- probability bounds
- KL divergence and cross entropy
- source coding
- logistic regression
- neural networks
- autoencoders, VAEs

Generative and Discriminative Models

posterior \propto likelihood \times prior

$$\hat{y} = \arg \max_y p(y|\mathbf{x}) \quad \text{MAP estimate}$$

$$\hat{y} = \arg \max_y p(\mathbf{x}|y) \cancel{p(y)} \quad \text{ML (maximum likelihood) estimate}$$

- we don't have $p(\mathbf{x}, y)$, instead we have $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$

- *generative* models

- we imagine some distribution $p(\mathbf{x}|y)$ that *generates* our data

example: imagine that $p(\mathbf{x}|y) \sim \mathcal{N}(\mathbf{u}_i, \Sigma_i)$

- *discriminative* models

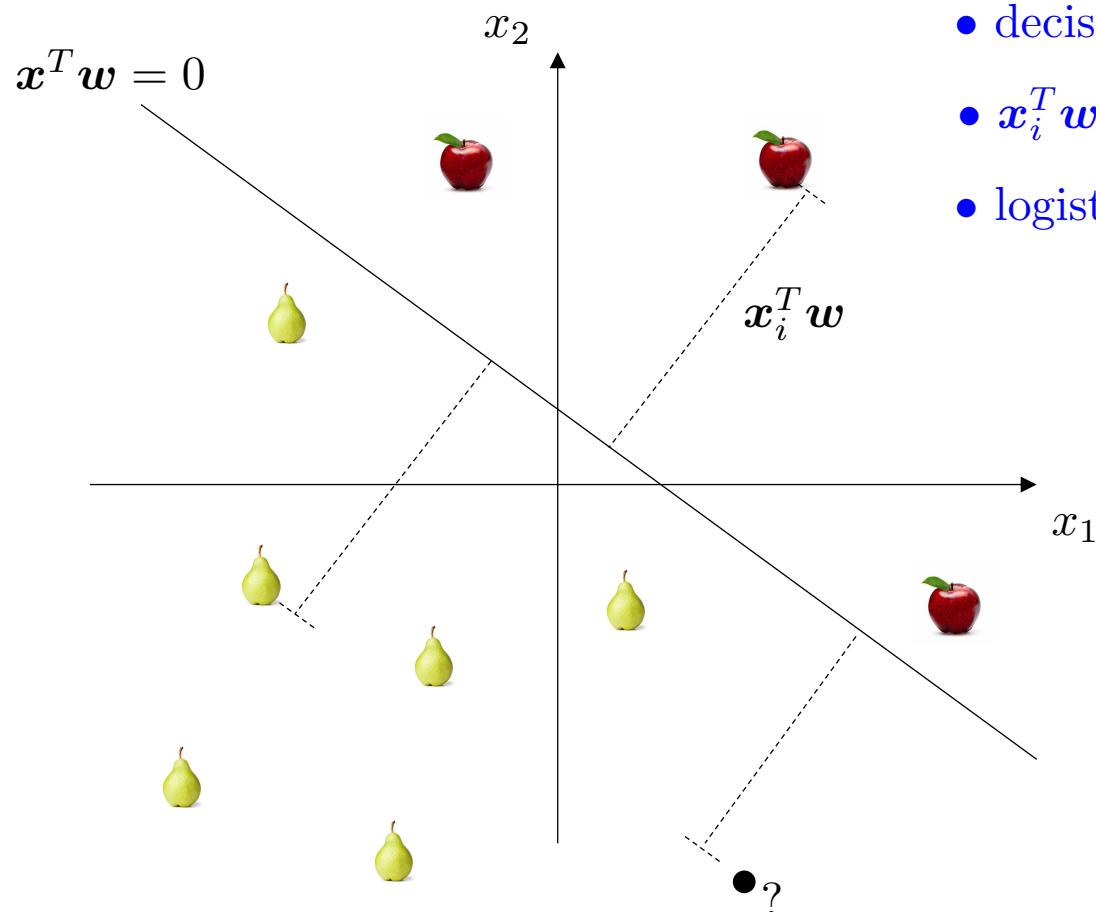
- discriminative classifiers directly find $p(y|\mathbf{x})$ to *discriminate* between classes.

directly fit $p(y|\mathbf{x})$ to data by treating \mathbf{x} as non-random

- we don't learn anything about the distribution of x

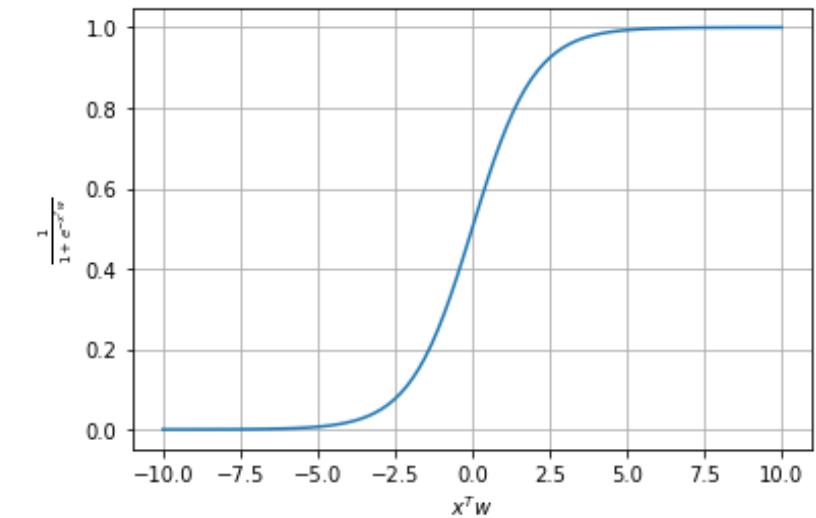
Two-Class Logistic Regression

- logistic regression is a *discriminative* classifier
directly fit $p(y|\mathbf{x})$ to data by treating \mathbf{x} as non-random



$$p(y=1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{x}^T \mathbf{w}}}$$

- decision boundary is $\mathbf{x}^T \mathbf{w} = 0$
- $\mathbf{x}_i^T \mathbf{w}$ is the signed distance from the decision boundary
- logistic function takes $\mathbf{x}_i^T \mathbf{w}$ and squishes it to make it a probability



- probability that ? is a pear/apple depends distance from boundary

Multi-Class Logistic Regression

- logistic regression is a *discriminative* classifier

$$y_i = \begin{cases} 1 & \text{with probability } \theta_i \\ -1 & \text{with probability } 1 - \theta_i \end{cases}$$

$$\theta_i = \frac{1}{1 + e^{-\mathbf{x}_i^T \mathbf{w}}} = \frac{e^{\mathbf{x}_i^T \mathbf{w}}}{e^{\mathbf{x}_i^T \mathbf{w}} + 1}$$

- multiclass (multinomial) logistic regression:

$$y_i = \begin{cases} 1 & \text{with probability } \theta_{i,1}(\mathbf{x}_i) \\ 2 & \text{with probability } \theta_{i,2}(\mathbf{x}_i) \\ \vdots & \\ c & \text{with probability } \theta_{i,C}(\mathbf{x}_i) \end{cases}$$

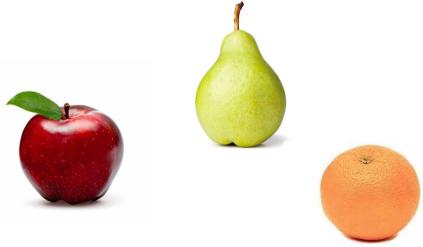
$$\theta_{i,c} = p(y = c | \mathbf{x}_i) = \frac{e^{\mathbf{x}_i^T \mathbf{w}_c}}{\sum_{c'=1}^C e^{\mathbf{x}_i^T \mathbf{w}_{c'}}$$

- softmax:

$$[\sigma(\mathbf{z})]_i = \frac{e^{z_i}}{\sum_j e^{z_j}} \quad \sigma(\mathbf{z}) = \frac{e^{\mathbf{z}}}{\sum_j e^{z_j}}$$

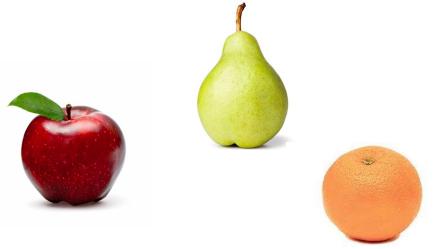
- softmax example:

$$\begin{matrix} [[-18.77687567] \\ [39.26910416] \\ [50.56209096] \\ [47.72731802] \\ [34.95432764]] \end{matrix} \xrightarrow{\hspace{10em}} \begin{matrix} [[0.] \\ [0.] \\ [0.9445] \\ [0.0555] \\ [0.]] \end{matrix}$$

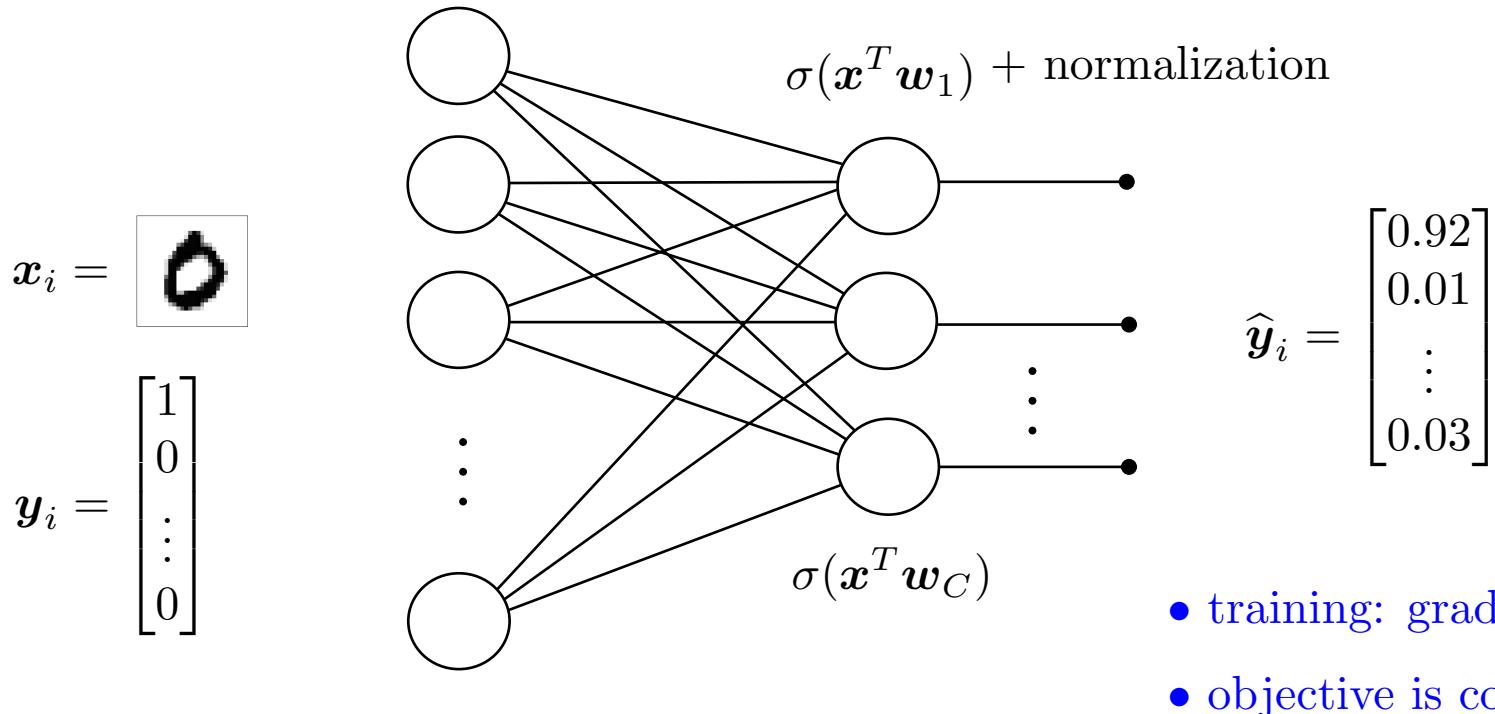


Multi-Class Logistic Regression, Single Layer

- logistic regression is a single *single layer of artificial neural network*



multiple neurons, single layer: $\hat{\mathbf{y}} = \sigma(\mathbf{W}^T \mathbf{x})$



- loss

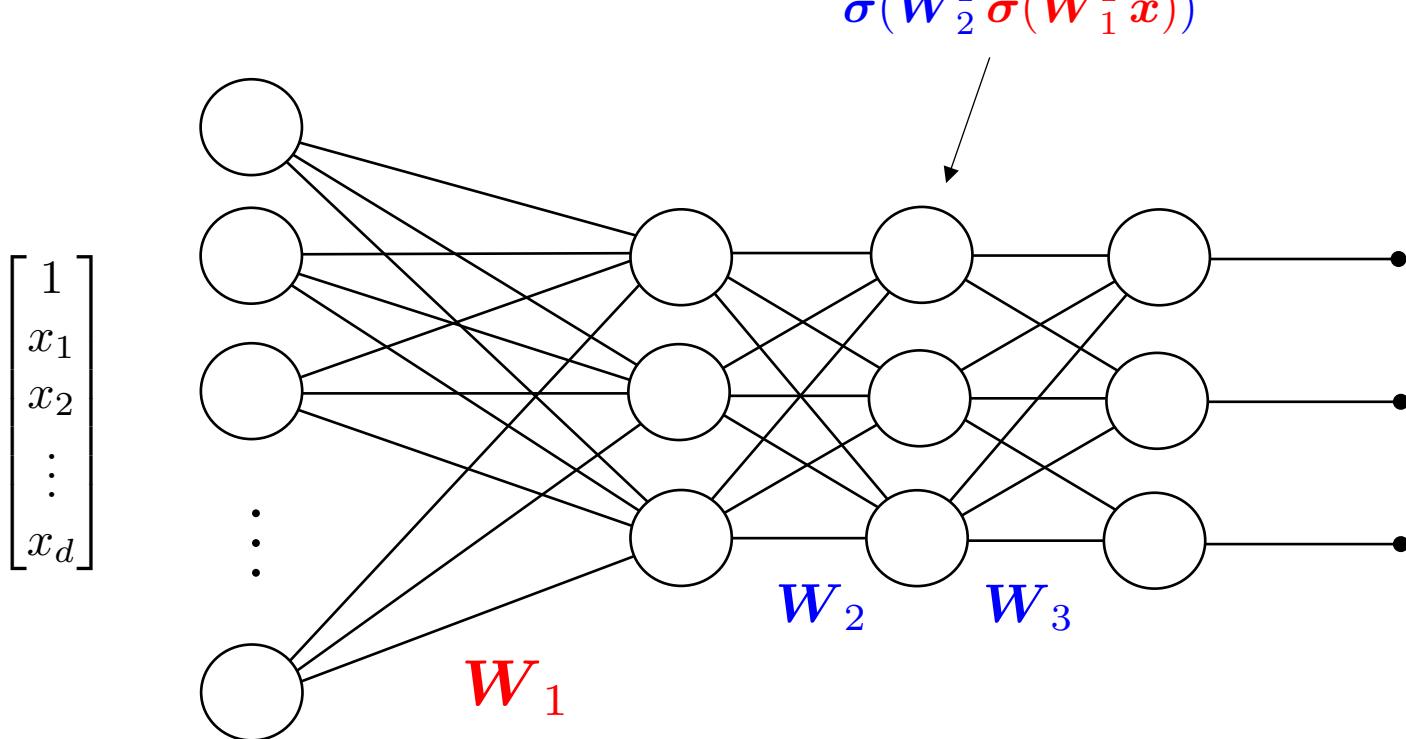
$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = D(\mathbf{y}_i || \hat{\mathbf{y}}_i)$$

$$\hat{\mathbf{y}}_i = \begin{bmatrix} 0.92 \\ 0.01 \\ \vdots \\ 0.03 \end{bmatrix}$$

- training: gradient descent/Newton's method
- objective is convex

Artificial Neural Networks and Logistic regression

- binary logistic regression:
 - single *artificial neuron*
- multiclass logistic regression:
 - single layer *artificial neural network*
- *feedforward* neural network:



single layer: $\hat{\mathbf{y}} = \sigma(\mathbf{W}^T \mathbf{x})$

two layers: $\hat{\mathbf{y}} = \sigma(\mathbf{W}_2^T \sigma(\mathbf{W}_1^T \mathbf{x}))$

ℓ -layers: $\hat{\mathbf{y}} = \sigma(\mathbf{W}_\ell^T \dots \sigma(\mathbf{W}_2^T \sigma(\mathbf{W}_1^T \mathbf{x})) \dots)$

- softmax:

$$[\sigma(\mathbf{z})]_i = \frac{e^z_i}{\sum_j e^{z_j}}$$

$$\sigma(\mathbf{z}) = \frac{e^{\mathbf{z}}}{\sum_j e^{z_j}}$$

- logistic (unnormalized):

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

KL Divergence, Entropy, Cross-Entropy

- entropy:

$$H(\mathbf{p}) = \sum_i p_i \log \frac{1}{p_i}$$

- number of bits required to encode a sequences of i.i.d. random variables

- cross-entropy:

$$H(\mathbf{p}, \mathbf{q}) = \sum_i p_i \log \frac{1}{q_i}$$

- number of bits required to encode a sequences of i.i.d. random variables from \mathbf{p} using a code optimized for \mathbf{q} .
- $H(\mathbf{p}, \mathbf{q})$ is big if $q_i \ll p_i$ for some i

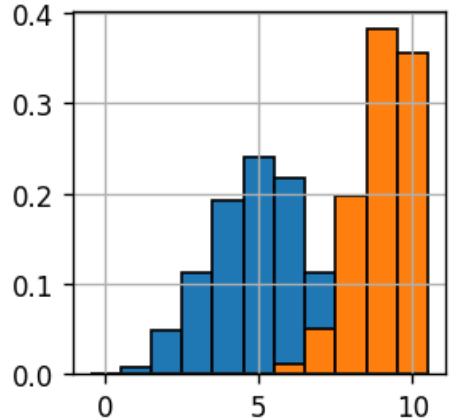
- Kullback-Leibler Divergence:

$$D(\mathbf{p}||\mathbf{q}) = \sum_i p_i \log \frac{p_i}{q_i}$$

- distance between \mathbf{p} and \mathbf{q}
- important in binary hypothesis testing
- KL divergence is non-negative: $D(p||q) \geq 0$ and $D(p||q) = 0$ if and only if $p(x) = q(x)$ for all x

- relationship:

$$H(\mathbf{p}, \mathbf{q}) = H(\mathbf{p}) + D(\mathbf{p}||\mathbf{q})$$



Source Coding

- Shannon's *source coding* theorem is fundamental to quantifying information
- applies to discrete data, i.e, discrete X
- interpretation:

n i.i.d. random variables each with entropy $H(X)$ can be compressed into $nH(X)$ bits with negligible risk of loss of information as n grows. Conversely, if compressed into fewer than $nH(X)$ bits, it is almost certain that information will be lost

[adapted from Cover and Thomas, 2006]

Probability Bounds

- imagine $X_i \stackrel{i.i.d.}{\sim} \mathcal{N}(\mu, \sigma^2)$
 - we're interested in the sample mean $\hat{\mu} = \frac{1}{n} \sum_i X_i$
- $$\hat{\mu} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

- how good is our estimate $\hat{\mu}$?

$$\begin{aligned}\mathbb{P}(|\hat{\mu} - \mu| \geq t) &= 2 \int_t^\infty \frac{1}{\sqrt{2\pi \frac{\sigma^2}{n}}} e^{-\frac{x^2}{2\sigma^2/n}} dx \\ &\leq 2e^{-t^2 n / \sigma^2}\end{aligned}$$

- *Hoeffding's Inequality.* Let X_1, X_2, \dots, X_n be i.i.d. bounded random variables with $X_i \in [a, b]$ and define $\hat{\mu} = \frac{1}{n} \sum X_i$, $\mu = E[X_i]$. Then

$$\mathbb{P}(|\hat{\mu} - \mu| \geq t) \leq 2e^{-\frac{2t^2 n}{(b-a)^2}}$$

- Hoeffding gives a tail bound that is basically like knowing the distributions are normal

Source Coding

$$X_i \stackrel{i.i.d.}{\sim} \text{bernoulli}(\theta)$$

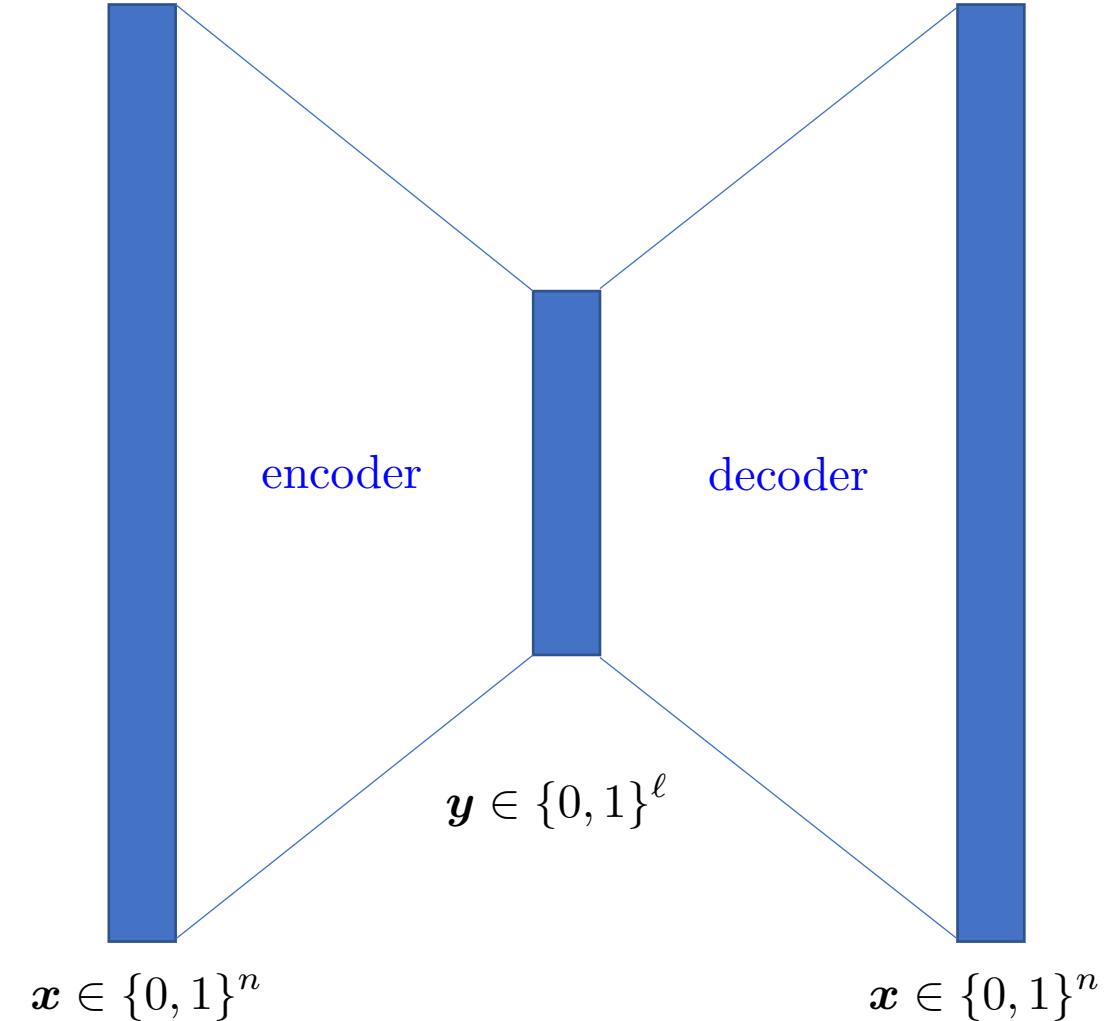
$$\boldsymbol{x} \sim \text{bernoulli}(\theta)$$

$$\mathbb{P}(|\hat{\theta} - \theta| \geq \delta) \leq 2e^{-2\delta^2 n}$$

$$\mathbb{P}(|k - n\theta| \geq m) \leq 2e^{-2m^2/n}$$

- only need to encode vectors that have between $n\theta - m$ and $n\theta + m$ ones.

$$\sum_{k=\theta n-m}^{\theta n+m} \binom{n}{k} \approx 2m \binom{n}{\theta n} \approx 2^{\log_2(2m)} \times 2^{nH(X)}$$



- compression ratio $\approx H(X)$
- crux of argument was showing $\binom{n}{\theta n} \approx 2^{nH(X)}$

Unsupervised Learning

- *supervised* machine learning is about **learning functions from data**

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$$

- *unsupervised* machine learning is about finding interesting patterns in data

$$\mathcal{D} = \{(x_i)\}_{i=1}^n$$

- *knowlege discovery*
- under-specified

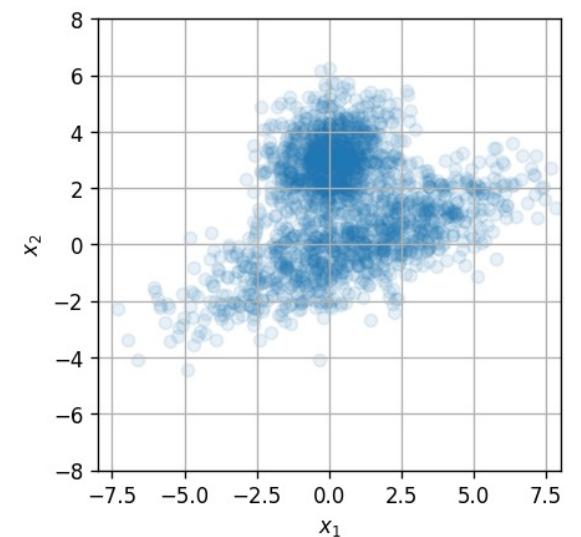
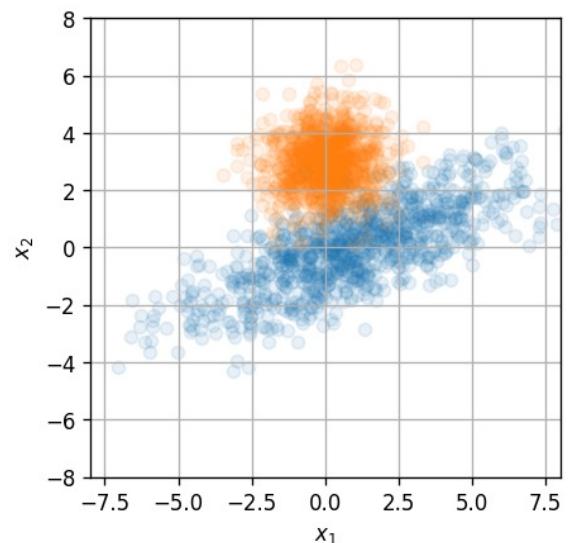
- learn $p(\mathbf{x})$ instead of $p(\mathbf{x}, y)$ or $p(y|\mathbf{x})$
- learn $p_\theta(\mathbf{x})$

- examples:

clustering

density estimation

community detection in graphs



Latent Variable Models

- goal: learn $p(\mathbf{x})$ from $\mathcal{D} = \{(\mathbf{x}_i)\}_{i=1}^n$

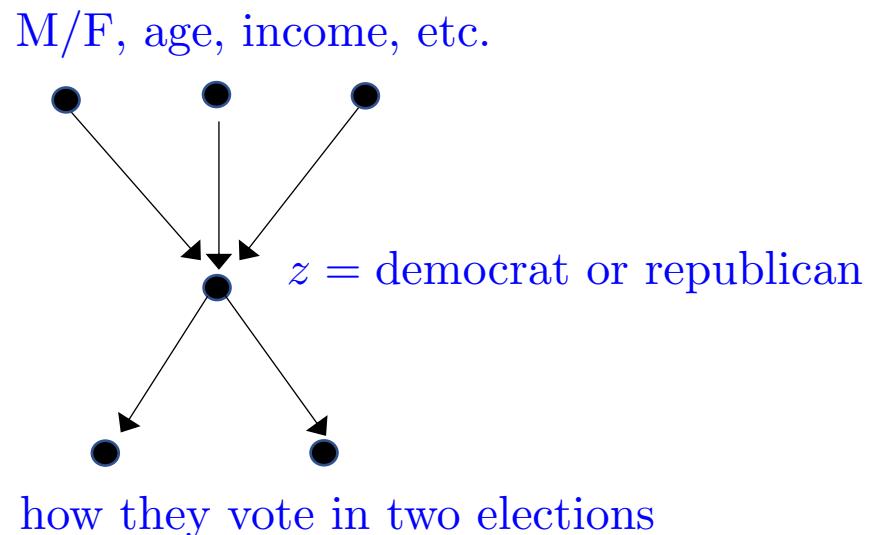
$$p(\mathbf{x}) = \int p(\mathbf{x}, z) dz$$

$$p(\mathbf{x}) = \int p(z) p(\mathbf{x}|z) dz$$

$$p(\mathbf{x}) = \sum_z p(\mathbf{x}, z)$$

$$p(\mathbf{x}) = \sum_z p(z) p(\mathbf{x}|z)$$

- z is a *latent* variable
- *might* be mathematically convenient
- maybe $p(\mathbf{x}|z)$ has a simple form
- z might be meaningful, or maybe just convenient
- MNIST: z = digit, pen width, slant, etc.
- *latent* means hidden or concealed
- if z is discrete, then *mixture* model



Expectation Maximization (EM) for GMMs

initialize means: $\boldsymbol{\mu}_k$, $j = 1, \dots, K$

initialize covariances: $\boldsymbol{\Sigma}_k = \mathbf{I}$

initialize priors: $\pi_k = \frac{1}{K}$

E-step: 1. compute the responsibility of each point to each cluster (*soft* assign each point):

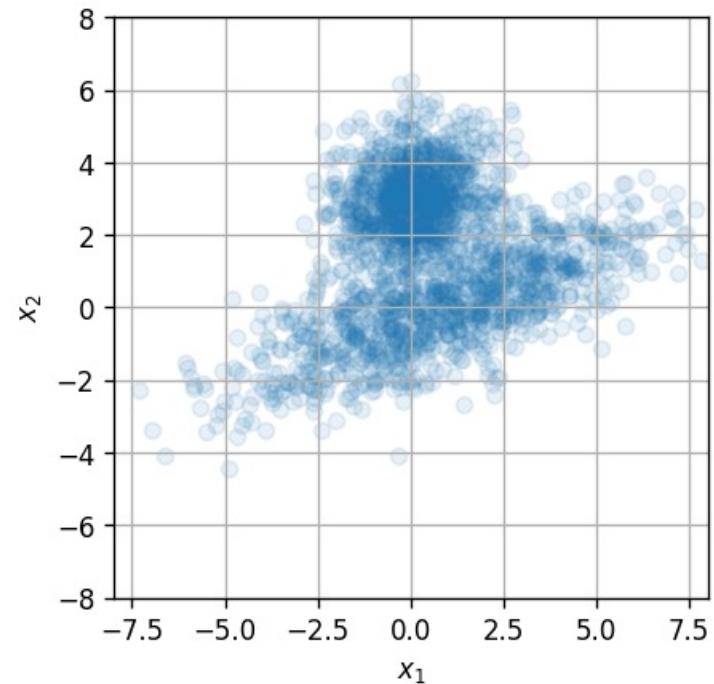
$$p(z = k | \mathbf{x}_i) = r_{ik} = \frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k'} \pi_{k'} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'})}$$

M-step: 2. estimate prior, mean and covariance for each mixture component

$$\boldsymbol{\mu}_k = \frac{1}{\sum_i r_{ik}} \sum_i r_{ik} \mathbf{x}_i \quad \boldsymbol{\Sigma}_k = \frac{1}{\sum_i r_{ik}} \sum_i r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T$$

$$\pi_k = \frac{1}{N} \sum_i r_{ik}$$

3. repeat 1, 2



Variational Inference

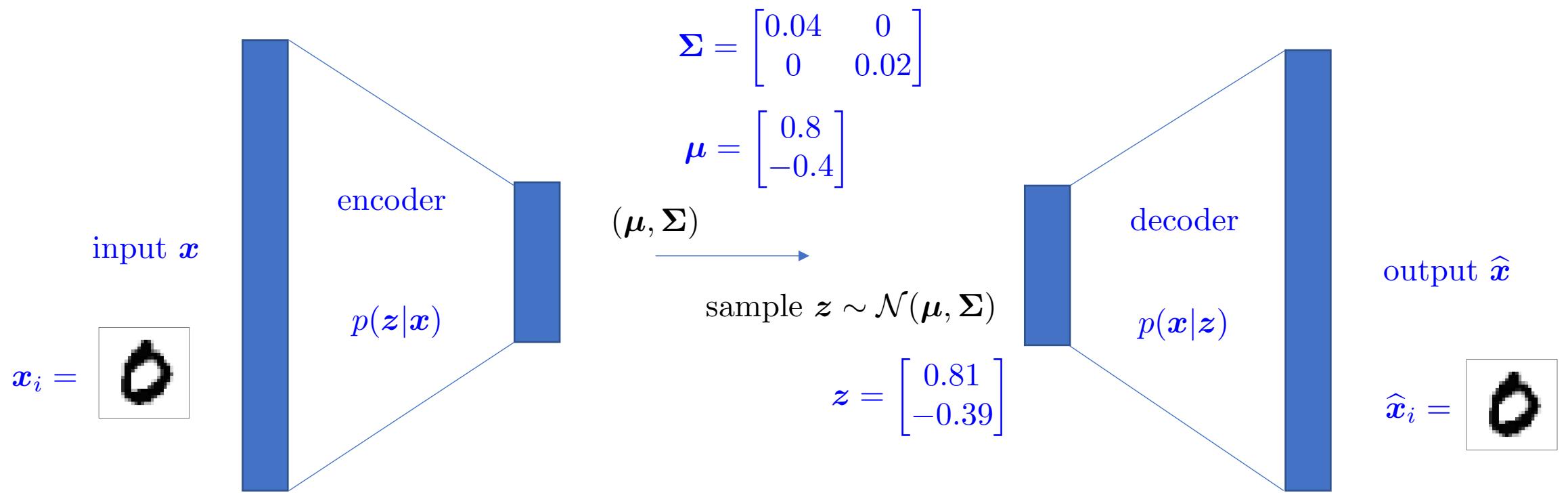
- goal: learn $p(\mathbf{z}|\mathbf{x})$

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})} = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}$$



- $p(\mathbf{z}|\mathbf{x})$ is the posterior over the hidden variable \mathbf{z}
- variational inference: estimate $p(\mathbf{z}|\mathbf{x})$ using some other distribution $q(\mathbf{z}|\mathbf{x})$
- restrict $q(\mathbf{z}|\mathbf{x})$ to be from some tractable family, i.e., $\mathcal{N}(\boldsymbol{\mu}(\mathbf{x}), \boldsymbol{\Sigma}(\mathbf{x}))$

Variational Autoencoders



$$\ell(x, \hat{x}) = ||\hat{x} - x||^2 + D(\mathcal{N}(\mu, \Sigma) || \mathcal{N}(0, I))$$

- ensure that $x \approx \hat{x}$
- try to force $p(z) \sim \mathcal{N}(0, I)$

Teaching Evaluations

<https://heliocampusac.wisc.edu/>