

**10417-617**  
**Deep Learning: Fall 2020**

Andrej Risteski

Machine Learning Department

**Lecture 10:**  
Intro to unsupervised  
learning

# Unsupervised learning

## ■ “Pure” Reinforcement Learning (cherry)

- ▶ The machine predicts a scalar reward given once in a while.
- ▶ **A few bits for some samples**

## ■ Supervised Learning (icing)

- ▶ The machine predicts a category or a few numbers for each input
- ▶ Predicting human-supplied data
- ▶ **10→10,000 bits per sample**

## ■ Unsupervised/Predictive Learning (cake)

- ▶ The machine predicts any part of its input for any observed part.
- ▶ Predicts future frames in videos
- ▶ **Millions of bits per sample**



■ (Yes, I know, this picture is slightly offensive to RL folks. But I'll make it up)

*Slide by Yann LeCun*

# Unsupervised learning

Learning from data **without** labels.

What can we hope to do:

**Task A:** Fit a parametrized **structure** (e.g. clustering, low-dimensional subspace, manifold) to data to reveal something meaningful about data. (**Structure learning**)

**Task B:** Learn a (parametrized) **distribution** *close* to data generating distribution. (**Distribution learning**)

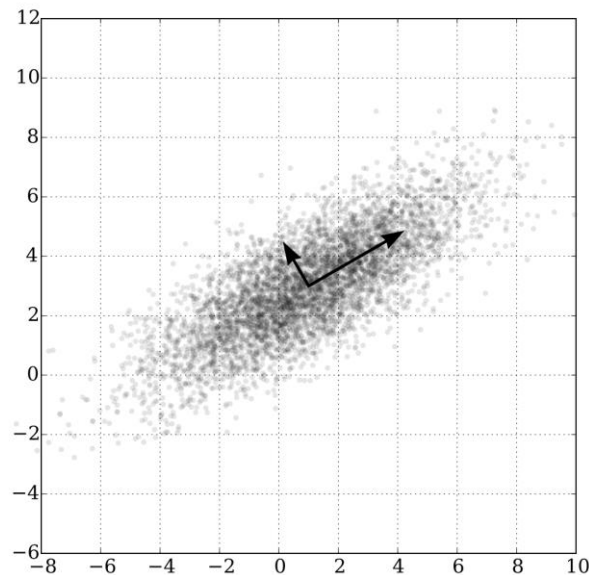
**Task C:** Learn a (parametrized) distribution that implicitly reveals an **“embedding”/“representation”** of data for downstream tasks. (**Representation/feature learning**)

*Entangled!* The “structure” and “distribution” often reveals an embedding.

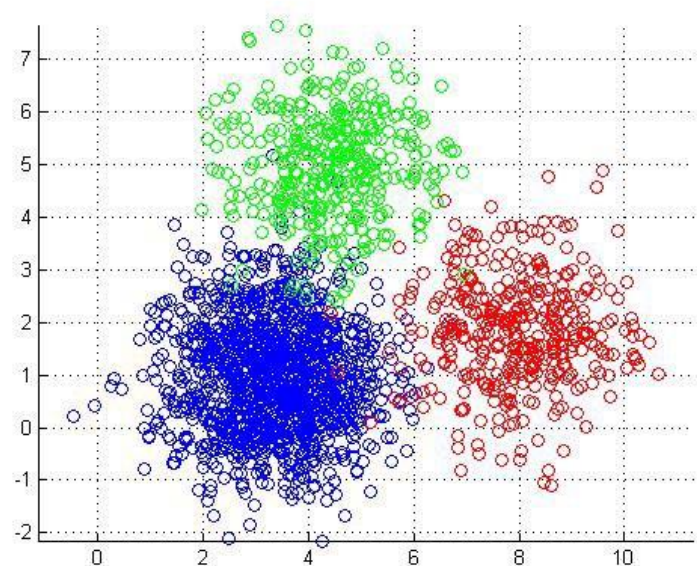
# Structure learning

Fit a parametrized **structure** (e.g. clustering, low-dimensional subspace) to data to reveal something meaningful about data.

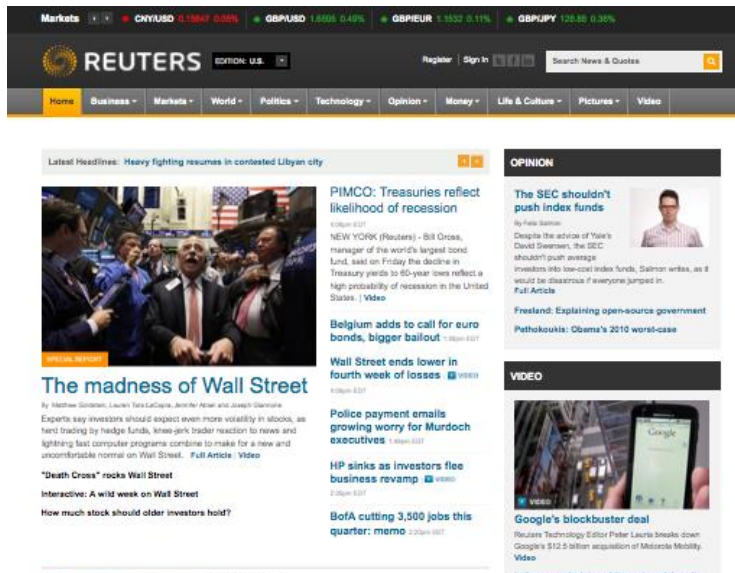
**PCA**(principal component analysis),  
direction of highest variance



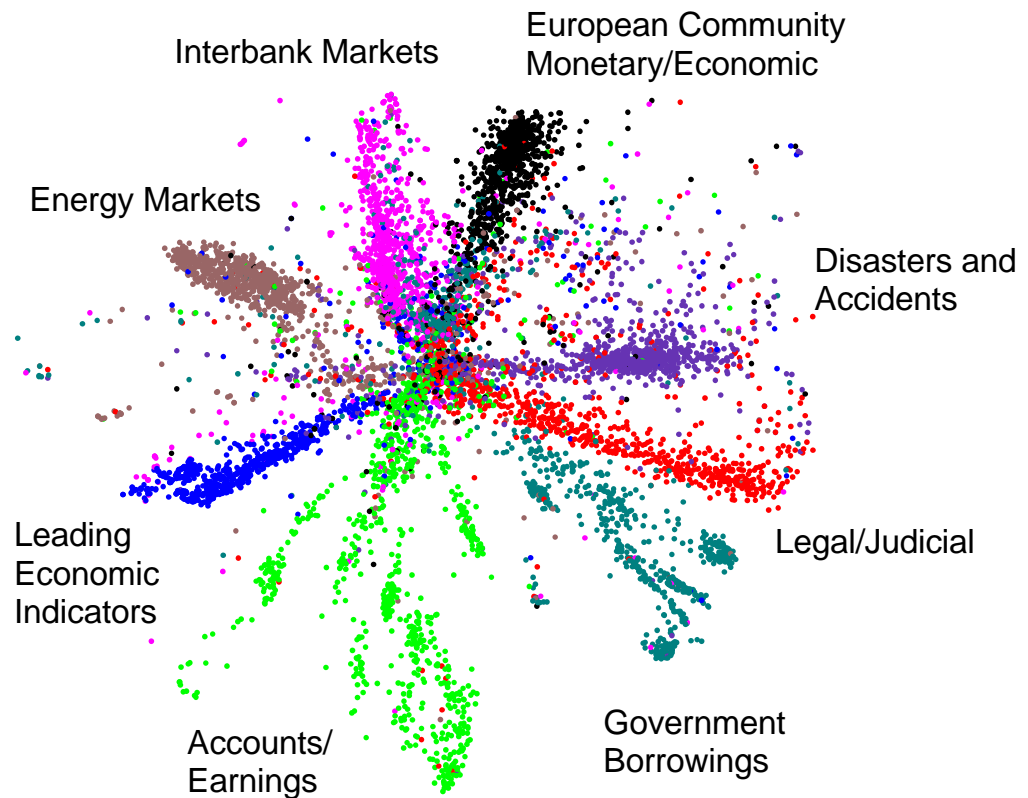
**Clustering**



# Structure learning



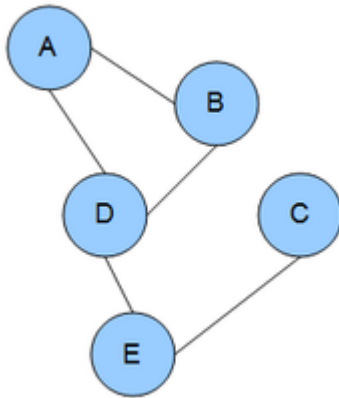
804,414 newswire stories



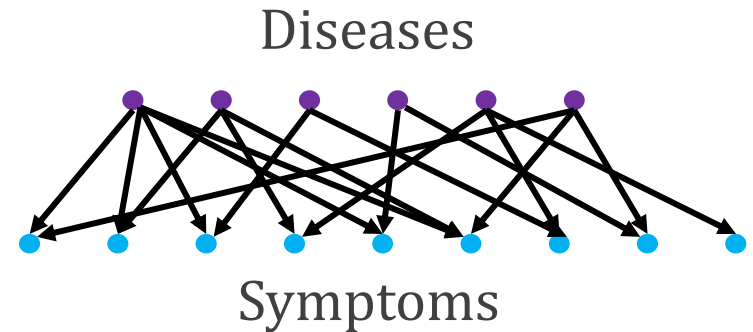
# Distribution learning

Some typical choices of parametrized distributions:

**Classical choices:** **fully-observed** graphical models (undirected and directed), **latent-variable** graphical models (mixture models, sparse coding, topic models).



**Markov Random Fields:**  
sparse independence  
structure: “A is independent of  
other vars, given B, D”

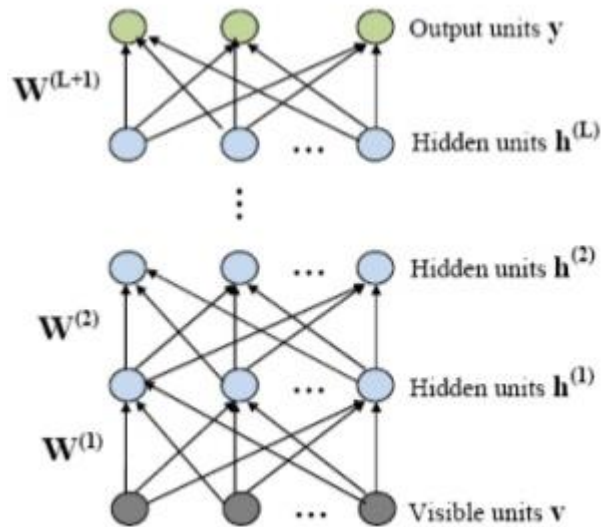


**Latent variable models:** data is  
“simple” conditioned on some  
unobserved (latent) variables

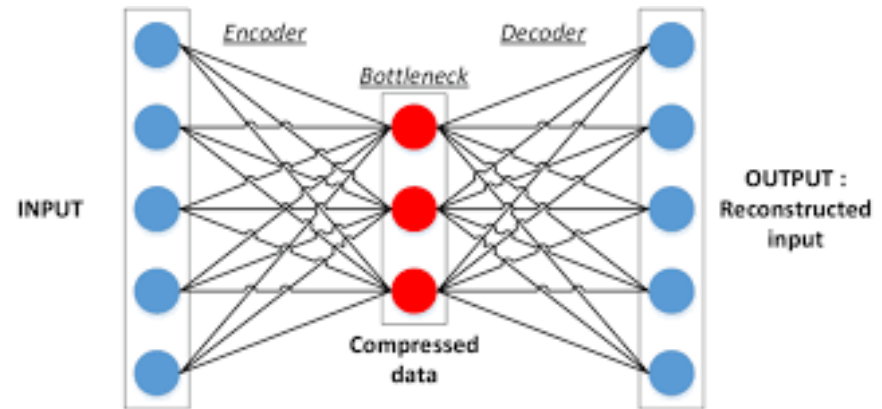
# Distribution learning

Some typical choices of parametrized distributions:

**Semi-modern choices:** deep Boltzmann machines, deep belief networks, (variational) auto-encoders, energy models.



**Deep Boltzmann machines, belief networks:** graphical model analogues of deep neural networks.



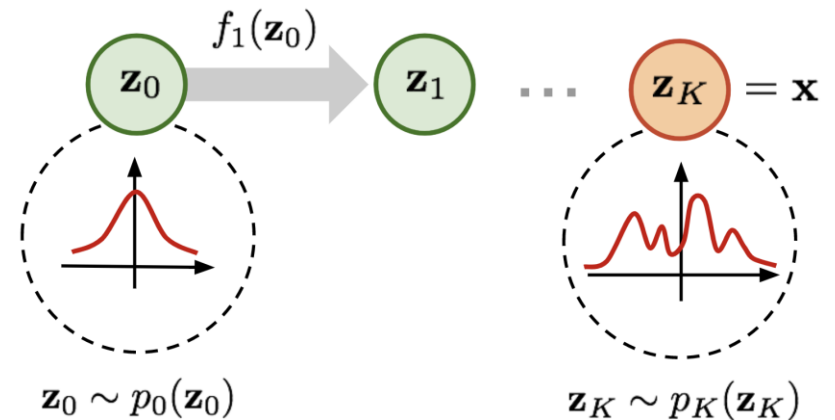
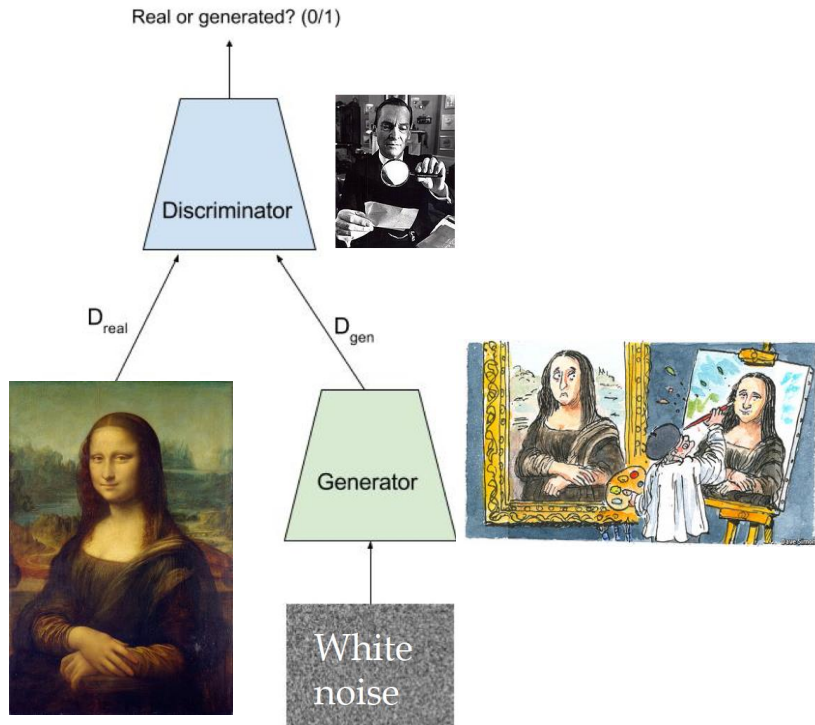
**(Variational) autoencoders:** model data by enforcing a latent space “bottleneck”:



# Distribution learning

Some typical choices of parametrized distributions:

**Modern choices:** generative adversarial networks, autoregressive models (pixelRNN, pixelCNN), flow models, etc.





# Distribution learning



Training  
Data(CelebA)



Model Samples (Karras et.al.,  
2018)

4 years of progression on Faces



2014



2015



2016



2017

Brundage et al.,  
2017

# Distribution learning



*Whichfaceisreal.com*

# Distribution learning



*BigGAN, Brock et al '18*



# Distribution learning

Conditional generative model  $P(\text{zebra images} | \text{horse images})$



Style Transfer



Input Image



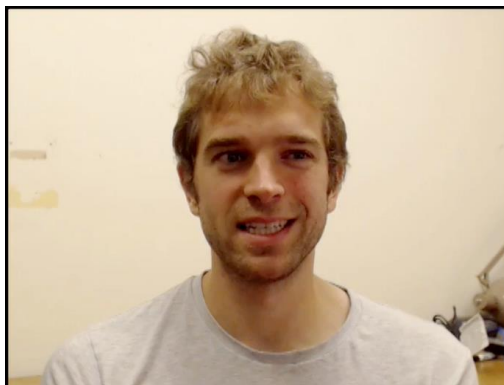
Monet



Van Gogh

# Distribution learning

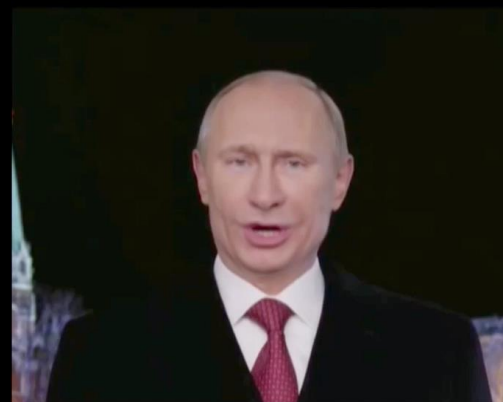
Source  
actor



Target  
actor



Real-time Reenactment



Real-time  
reenactment

Reenactment Result

# Representation learning and self-supervised learning

Given **unlabeled** data, **design supervised tasks** that induce a good representation for downstream tasks.

No good mathematical formalization, but the intuition is to “force” the predictor used in the task to learn something “**semantically meaningful**” about the data.

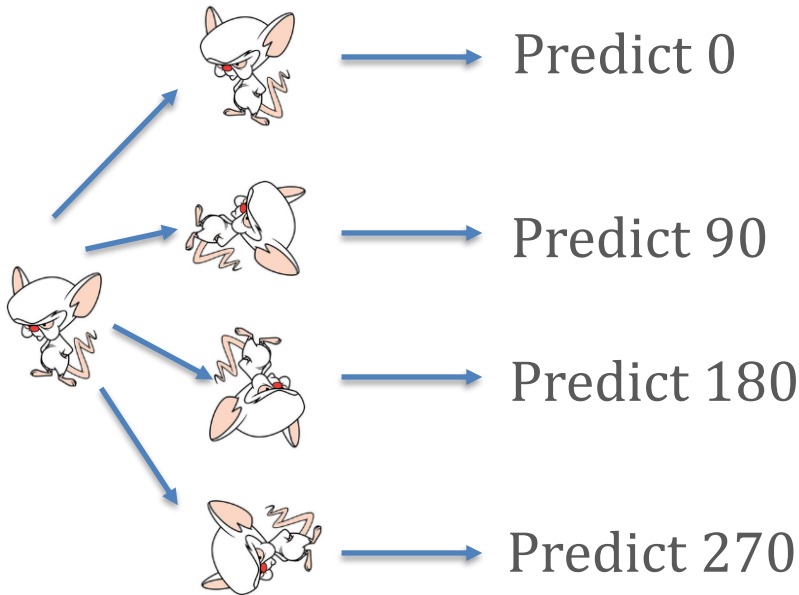
**Examples in NLP:** predict next word, given previous 5 words; predict middle word, given surrounding 5 words; etc.



# Representation learning and self-supervised learning

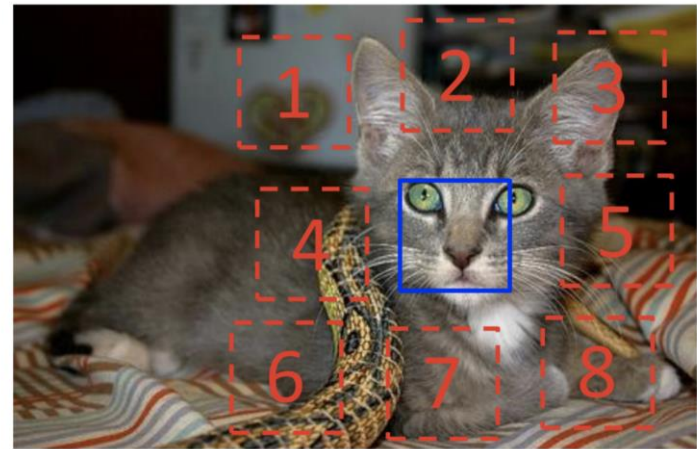
**Examples in vision:** a lot, and quite different in nature.

## Rotation prediction



Predict one of four angles  
an image is rotated by

## Jigsaw puzzles



$$X = \left( \begin{array}{c} \text{cat face} \\ \text{cat ear} \end{array} \right); Y = 3$$

Predict position of second  
piece wrt to first

# Relationships between the tasks

## **Structure learning and feature learning often blend**

E.g. low dimensional features in PCA or cluster a point belongs to in clustering can be viewed as features.

## **Structure learning/feature learning is in general weaker than distribution learning:**

E.g. methods like PCA/clustering can't be used to generate new samples.

*Feature, not a bug*: it has been argued that self-supervised learning works because the task we are solving is **easier** (both computationally and statistically)

# Relationships between the tasks

## Distribution learning often implies representation learning:

Many distributions we fit are **latent-variable** models (i.e. model the joint distribution between some latent variables  $\mathbf{h}$  and the observed data  $\mathbf{x}$ )

$$P_{\theta}(x, h) = P_{\theta}(h)P_{\theta}(x|h)$$


The latent variables are often viewed as a “**representation**”.

The **posterior** distribution  $P_{\theta}(h|x)$  captures a distribution over representations, given some values of the observed data.

However, a-priori, distribution  $P_{\theta}(h|x)$  is a-priori not an easy distribution to **approximate/sample** from!

$$P_{\theta}(h|x) = \frac{P_{\theta}(x|h)p(h)}{\int_{h'} p(h')p(x|h')}$$

*Hard high-dimensional  
sum/integral*



# Training techniques: Likelihood-based vs likelihood-free

**Two typical** families of training algorithms:

**Likelihood-based:** maximize the likelihood of the data under the model (possibly with some approximations)

$$\max_{\theta} \sum_{\text{samples } x_i} \log p_{\theta}(x_i)$$

In the limit of infinite number of samples:

$$\begin{aligned} \max_{\theta} \mathbb{E}_{x \sim p} \log p_{\theta}(x) &= -(\mathbb{E}_{x \sim p} \log \frac{1}{p_{\theta}(x)} + \mathbb{E}_{x \sim p} \log p(x) - \mathbb{E}_{x \sim p} \log p(x)) \\ &= -( \quad \quad \quad KL(p || p_{\theta}) \quad \quad \quad + \quad H(p) ) \end{aligned}$$

Hence, we are minimizing KL divergence ( $H(p)$  is a constant).

# Training techniques: Likelihood-based vs likelihood-free

**Two typical** families of training algorithms:

**Likelihood-based:** maximize the likelihood of the data under the model (possibly with some approximations)

## *Pros*

**Easy training:** can just maximize via gradient descent.

**Evaluation:** evaluating the fit of the model can be done by evaluating the likelihood (on test data)

## *Cons*

**Larger models needed:** likelihood objective is hard, to fit well need very big model

**Likelihood encourages averaging:** produced samples tend to be blurrier, as likelihood encourages “coverage” of training data.

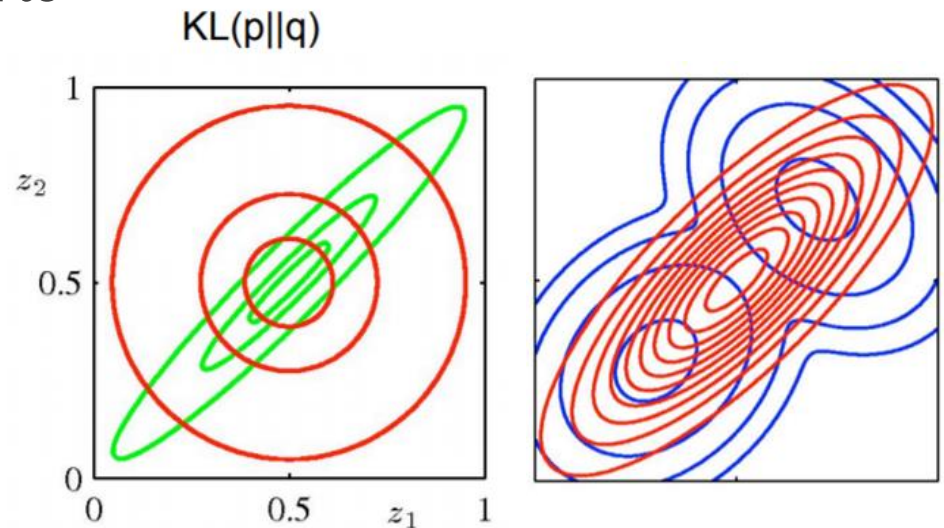
# Training techniques: Likelihood-based vs likelihood-free

$$\text{KL}(p||q) = - \int p(\mathbf{Z}) \ln \frac{q(\mathbf{Z})}{p(\mathbf{Z})} d\mathbf{Z}.$$

There is a large positive contribution to the KL divergence from regions of  $\mathbf{Z}$  space in which:

- $q(\mathbf{Z})$  is near zero,
- unless  $p(\mathbf{Z})$  is also close to zero.

Minimizing  $\text{KL}(p||q)$  leads to distributions  $q(\mathbf{Z})$  that **are nonzero in regions where  $p(\mathbf{Z})$  is nonzero.**





# Training techniques: Likelihood-based vs likelihood-free



Faces generated using a trained VAE

# Training techniques: Likelihood-based vs likelihood-free

**Two typical** families of training algorithms:

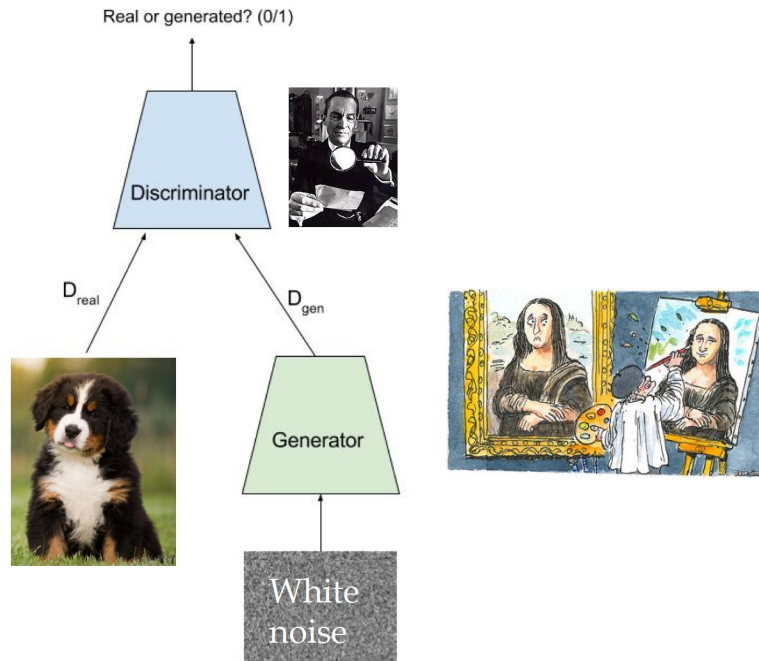
**Likelihood-based:** maximize the likelihood of the data under the model (possibly with some approximations)

*Typical approximations used:* variational inference (optimize tractable deterministic approximation of posteriors), MCMC methods (idea: approximate difficult quantities like posteriors with sampling)

**Likelihood-free:** use a surrogate loss – e.g. in GANs, train a discriminator to tell real and generated samples apart; noise-contrastive training: encourage model to put probability mass away from “fake” samples.

# Training techniques: Likelihood-based vs likelihood-free

**Likelihood-free:** use a **surrogate loss** – e.g. in GANs, train a discriminator to tell real and generated samples apart; **noise-contrastive training:** encourage model to put probability mass away from “fake” samples.



# Training techniques: Likelihood-based vs likelihood-free

**Likelihood-free:** use a **surrogate loss** – e.g. in GANs, train a discriminator to tell real and generated samples apart; **noise-contrastive training:** encourage model to put probability mass away from “fake” samples.

## *Pros*

**Better objective, smaller models needed:** objective itself (i.e. discriminator) is “learned” – can result in visually better images w/ smaller models.

## *Cons*

**Unstable training:** typically min-max (saddle point) problems.

**Evaluation:** no way to evaluate likelihood, so no way to evaluate quality of fit.

# The classics: PCA

Figure 1: Population structure within Europe.

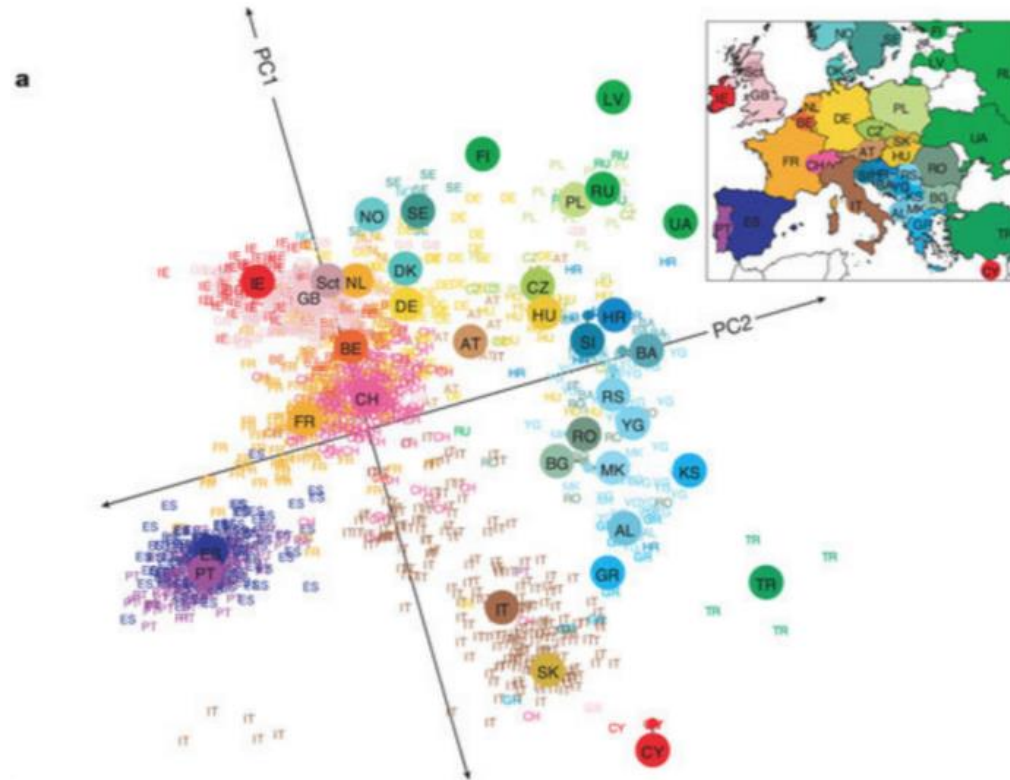


Figure 6: Plot from [1], depicting genomes for 1387 Europeans projected onto top 2 principal components. Colors/labels of datapoints correspond to geographic location of the individuals. Map of Europe (with same coloring) included in upper right for reference.

# The classics: PCA

**Goal:** find a  $k$ -dimensional (linear) subspace explaining most of the variance in the data.

Assume the data is centered, that is  $\mathbb{E}[x] = 0$

*Warmup:* let  $k=1$ .

$$\operatorname{argmax}_{\{v: \|v\|=1\}} \frac{1}{m} \sum_{\text{samples } x_i} \langle v, x_i \rangle^2$$




As  $m \rightarrow \infty$ , we get  $\mathbb{E}[\langle v, x \rangle^2]$ ,  
i.e. variance.



# The classics: PCA

**Goal:** find a  $k$ -dimensional (linear) subspace explaining most of the variance in the data.

$$\operatorname{argmax}_{\{v_1, v_2, \dots, v_k\}} \frac{1}{m} \sum_{\text{samples } x_i} (\text{length of } x_i \text{ on } \operatorname{span}(v_1, v_2, \dots, v_k))^2$$

$$= \operatorname{argmax}_{\text{orthonormal } \{v_1, v_2, \dots, v_k\}} \frac{1}{m} \sum_{\text{samples } x_i} \sum_{j=1}^k \langle x_i, v_j \rangle^2$$


As  $m \rightarrow \infty$ , we get  $\mathbb{E}[\sum_j \langle v_j, x \rangle^2]$

# The classics: PCA

*How to do this efficiently?*

$$\operatorname{argmax}_{\text{orthonormal } \{v_1, v_2, \dots, v_k\}} \frac{1}{m} \sum_{\text{samples } x_i} \sum_{j=1}^k \langle x_i, v_j \rangle^2$$

A convenient rewrite:

$$= \operatorname{argmax}_{\text{orthonormal } \{v_1, v_2, \dots, v_k\}} \frac{1}{m} \sum_{\text{samples } x_i} \sum_{j=1}^k (v_j^T x_i) (x_i^T v_j)$$

$$= \operatorname{argmax}_{\text{orthonormal } \{v_1, v_2, \dots, v_k\}} \frac{1}{m} \sum_{\text{samples } x_i} \sum_{j=1}^k v_j^T (x_i x_i^T) v_j$$

$$= \operatorname{argmax}_{\text{orthonormal } \{v_1, v_2, \dots, v_k\}} \sum_{j=1}^k v_j^T \underbrace{\left( \frac{1}{m} \sum_{\text{samples } x_i} (x_i x_i^T) \right)}_{\text{:= D, (covariance matrix)}} v_j$$

:= D, (covariance matrix)

# The classics: PCA

How to do this efficiently? – **Singular Value Decomposition!!**

$$\operatorname{argmax}_{\text{orthonormal } \{v_1, v_2, \dots, v_k\}} \frac{1}{m} \sum_{\text{samples } x_i} \sum_{j=1}^k \langle x_i, v_j \rangle^2$$

A convenient rewrite:

$$\operatorname{argmax}_{\text{orthonormal } \{v_1, v_2, \dots, v_k\}} \sum_{j=1}^k v_j^T D v_j$$

If  $D$  is **diagonal** (w/ positive entries), and we sort the entries s.t.  $D_{11} \geq D_{22} \dots \geq D_{dd}$ , it's easy to see max is  $\sum_{j=1}^k D_{jj}$ . Namely:

$$\sum_{j=1}^k v_j^T D v_j = \sum_j \sum_i (v_j)_i^2 D_{ii} \leq \sum_{j=1}^k D_{jj}, \text{ as } (v_j)_i^2 = 1$$

The corresponding argmax is  $e_1, e_2, \dots, e_k$ .

# The classics: PCA

How to do this efficiently? – **Singular Value Decomposition!!**

$$\operatorname{argmax}_{\text{orthonormal } \{v_1, v_2, \dots, v_k\}} \frac{1}{m} \sum_{\text{samples } x_i} \sum_{j=1}^k \langle x_i, v_j \rangle^2$$

A convenient rewrite:

$$\operatorname{argmax}_{\text{orthonormal } \{v_1, v_2, \dots, v_k\}} \sum_{j=1}^k v_j^T D v_j$$

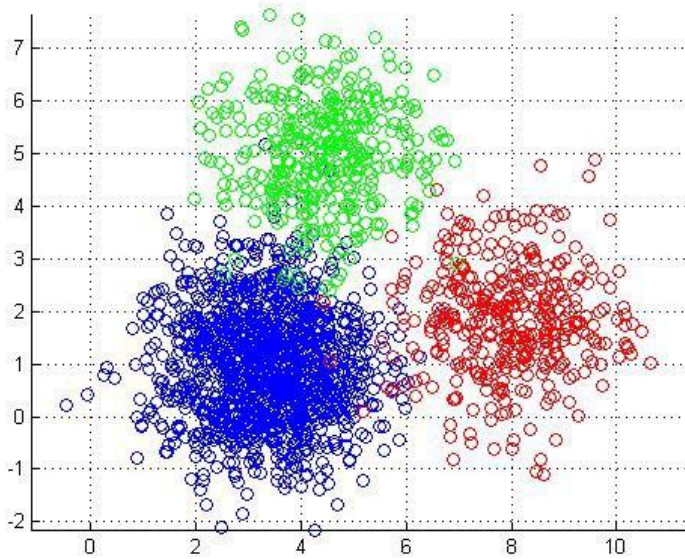
If  $D$  is not diagonal, write  $D = U \tilde{D} U^T$  to reduce to diagonal case:

$$\operatorname{argmax}_{\text{orthonormal } \{v_1, v_2, \dots, v_k\}} \sum_{j=1}^k v_j^T (U \tilde{D} U^T) v_j = \operatorname{argmax}_{\text{orthonormal } \{v_1, v_2, \dots, v_k\}} \sum_{j=1}^k (U^T v_j)^T \tilde{D} (U^T v_j)$$

Since  $U^T$  is orthogonal,  $\{U^T v_1, U^T v_2, \dots, U^T v_k\}$  also are orthonormal!

So, the max is  $\sum_{j=1}^k \tilde{D}_{jj} = \sum_{j=1}^k \lambda_j(D)$  and the argmax are the **top k eigenvectors** of  $D$

# The classics: clustering

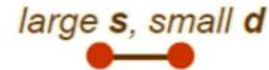
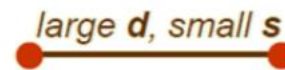


**Goal:** group the data into clusters of nearby points.

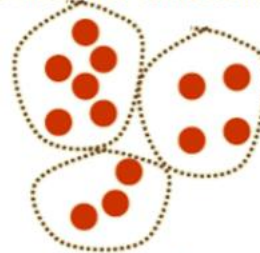
## What's needed for clustering?

### 1. Proximity measure, *either*

- similarity measure  $s(x_i, x_k)$ : large if  $x_i, x_k$  are similar
- dissimilarity (or distance) measure  $d(x_i, x_k)$ : small if  $x_i, x_k$  are similar



### 2. Criterion function to evaluate a clustering



*good clustering*



*bad clustering*

### 3. Algorithm to compute clustering

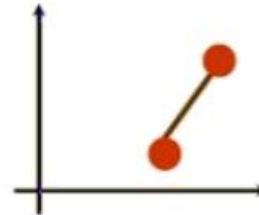
# The classics: clustering

## Popular distance metrics:

- Euclidean distance

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{\sum_{k=1}^d (\mathbf{x}_i^{(k)} - \mathbf{x}_j^{(k)})^2}$$

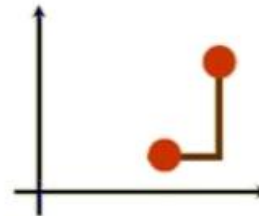
- translation invariant



- Manhattan (city block) distance

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sum_{k=1}^d |\mathbf{x}_i^{(k)} - \mathbf{x}_j^{(k)}|$$

- approximation to Euclidean distance, cheaper to compute



- They are special cases of **Minkowski distance**:

$$d_p(\mathbf{x}_i, \mathbf{x}_j) = \left( \sum_{k=1}^m |\mathbf{x}_{ik} - \mathbf{x}_{jk}|^p \right)^{\frac{1}{p}}$$

(p is a positive integer)



# The classics: clustering

## Criterion functions:

### Intra-cluster cohesion

- Cohesion measures how near the data points in a cluster are to the cluster “center”.

### Inter-cluster separation

- Separation means that different cluster centroids should be far away from one another.

In most applications, expert judgments are still the key

# The classics: clustering

How many clusters?



- Possible approaches
  1. fix the number of clusters to  $k$
  2. find the best clustering according to the criterion function (number of clusters may vary)

# K-means clustering

If the distance metric is the **Euclidean distance**, and the measure of cohesion is the **average distance from the centroid**: we get the **k-means objective**.

$$\operatorname{argmin}_{\{r_{nk}, \mu_k\}} \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \mu_k\|^2$$

*Is point  $n$  in cluster  $k$ ?*

*Centroid of  $k$ -th cluster*

# K-means clustering

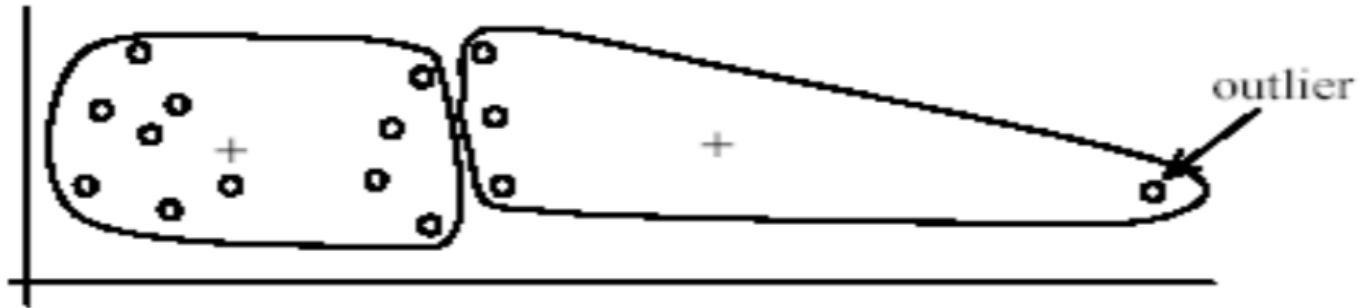
A natural iterative algorithm, which as we will see later is a variant of the **EM (expectation-maximization)** algorithm:

```
Input: Data set  $X = \{x^{(1)}, x^{(2)}, \dots, x^{(m)} | x^{(i)} \in \mathbb{R}^n\}$ 
Output: Cluster centroids  $\mu_{i=1, \dots, k} \in \mathbb{R}^n$ ; Cluster assignments  $c \in \mathbb{Z}$ 
1 Initialize  $k$  cluster centroids  $\mu_1, \dots, \mu_k \in \mathbb{R}^n$  randomly from  $X$ ;
2 repeat
3   for  $i = 1, \dots, m$  do // Update cluster assignments
4     | set  $c^{(i)} = \arg \min_j \|x^{(i)} - \mu_j\|^2$ ;
5   end
6   for  $j = 1, \dots, k$  do // Update cluster centroids
7     | set  $\mu_j = \frac{\sum_{i=1}^m 1_{\{c^{(i)}=j\}} x^{(i)}}{\sum_{i=1}^m 1_{\{c^{(i)}=j\}}}$ ;
8   end
9 until Convergence;
10 return  $\mu$  and  $c$ ;
```

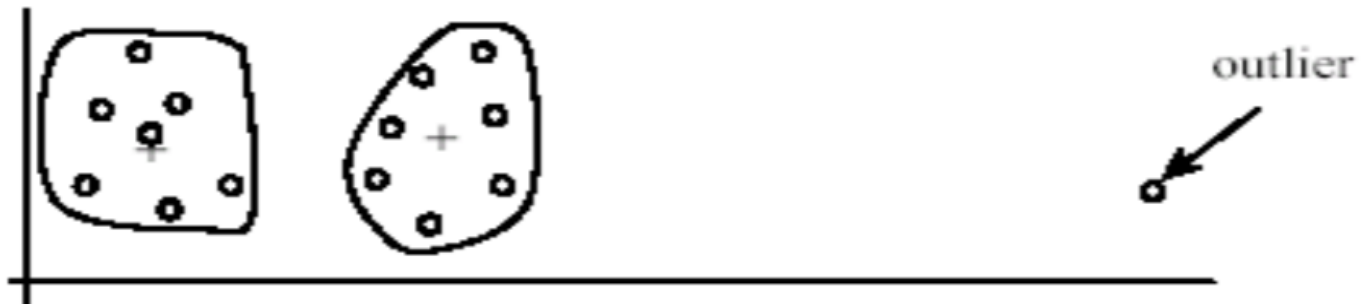
Algorithm 1: Algorithm of batch-version for K-means

# Some weaknesses

Very sensitive to outliers:



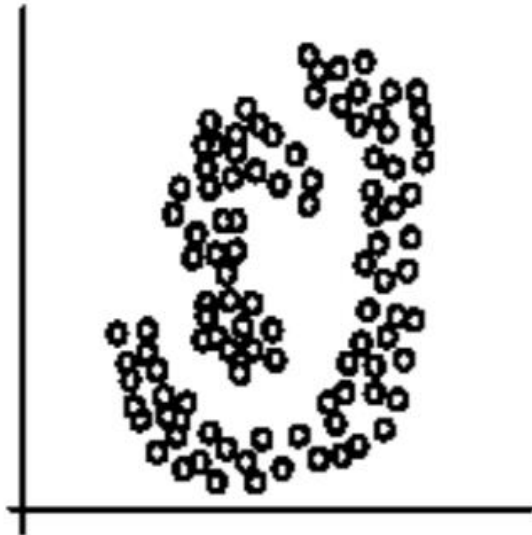
(A): Undesirable clusters



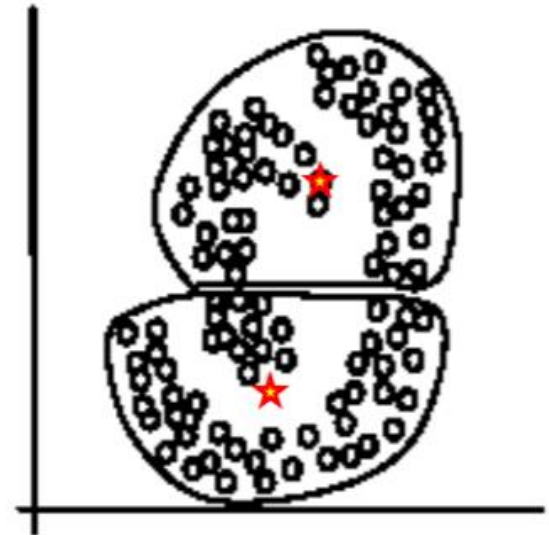
(B): Ideal clusters

# Some weaknesses

Not suitable for non-spherical clusters:



(A): Two natural clusters



(B):  $k$ -means clusters