

# Visualisation and Dimension Reduction

Christophe Ambroise 

*christophe.ambroise@univ-evry.fr*

*UEVE, UMR CNRS 8071*

February 20, 2025

1

## Evaluation

A Natural Language Processing (NLP) project to be completed in pairs:

- Code (in a Python notebook or R Markdown)
- Presentation (scheduled for April 4, 2025)

2

# Factor Analyser

3

## Factor Analysis

- Using discrete latent variables provides limited summary (clustering)
- An alternative is to use a vector of real-valued latent variables,  $\mathbf{z} \in \mathbb{R}^L$ .
- “Factor analysis (FA) is a statistical method used to describe variability among observed, correlated variables in terms of a potentially lower number of unobserved variables called factors.” Wikipedia quote.
- PCA and FA are related, but not identical.

4

# The model of factor analysis

We consider the observation  $\mathbf{x} \in \mathbb{R}^D$

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

where

- the noise  $\boldsymbol{\epsilon} \sim \mathcal{N}_D(\mathbf{0}, \boldsymbol{\Psi})$
- the hidden (latent) vector  $\mathbf{z} \sim \mathcal{N}_L(\mathbf{0}, \mathbf{I}_L)$

$$p(\mathbf{x}|\mathbf{z}, \theta) = \mathcal{N}(\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \boldsymbol{\Psi})$$

the mean is a linear function of the (hidden) inputs

- $\mathbf{W}$  is a  $D \times L$  matrix, known as the factor loading matrix,
- $\boldsymbol{\Psi}$  is a  $D \times D$  covariance matrix that we take to be diagonal

The special case in which  $\boldsymbol{\Psi} = \sigma^2 \mathbf{I}$  is called probabilistic principal components analysis or PPCA.

5

## Reminder: Joint and conditional Gaussian distribution (see Murphy chapter 4)

Let us recall that if

- $\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}_z, \boldsymbol{\Sigma}_{zz})$  and  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_{xx})$  and  $\text{cov}(\mathbf{z}, \mathbf{x}) = \boldsymbol{\Sigma}_{zx}$

then

$$p(\mathbf{z}, \mathbf{x}) = N\left(\begin{bmatrix} \mathbf{z} \\ \mathbf{x} \end{bmatrix} \mid \begin{bmatrix} \boldsymbol{\mu}_z \\ \boldsymbol{\mu}_x \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{zz} & \boldsymbol{\Sigma}_{zx} \\ \boldsymbol{\Sigma}_{xz} & \boldsymbol{\Sigma}_{xx} \end{bmatrix}\right)$$

and

$$p(\mathbf{z}, \mathbf{x}) = p(\mathbf{z}|\mathbf{x})p(\mathbf{x}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{z|x}, \boldsymbol{\Sigma}_{z|x})\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_{xx})$$

where

- $\boldsymbol{\mu}_{z|x} = \boldsymbol{\mu}_z + \boldsymbol{\Sigma}_{zx} \boldsymbol{\Sigma}_{xx}^{-1} (\boldsymbol{x} - \boldsymbol{\mu}_x)$
- $\boldsymbol{\Sigma}_{z|x} = \boldsymbol{\Sigma}_{zz} - \boldsymbol{\Sigma}_{zx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xz}$

## Marginal and posterior distribution

### Marginal distribution

$$\boldsymbol{x} \sim \mathcal{N}_D(\boldsymbol{\mu}, \boldsymbol{\Sigma}_{xx} = \boldsymbol{W}\boldsymbol{W}^T + \boldsymbol{\Psi})$$

### Posterior distribution

$$\boldsymbol{z}|\boldsymbol{x} \sim \mathcal{N}_L(\boldsymbol{\mu}_{z|x}, \boldsymbol{\Sigma}_{z|x})$$

where

- $\boldsymbol{\Sigma}_{z|x} = (\boldsymbol{I}_L + \boldsymbol{W}^T \boldsymbol{\Psi}^{-1} \boldsymbol{W})^{-1} = \boldsymbol{S}$
- $\boldsymbol{\mu}_{z|x} = \boldsymbol{\Sigma}_{z|x} \boldsymbol{\Sigma}_{xx}^{-1} (\boldsymbol{x} - \boldsymbol{\mu}) = \boldsymbol{S} \boldsymbol{W}^T \boldsymbol{\Psi}^{-1} (\boldsymbol{x} - \boldsymbol{\mu})$

### Exercise

Demonstrate the above formulas

# Estimation

## The mean $\mu$

can be estimated by maximum likelihood

$$\mu_{mle} = \bar{x}$$

## $W$ and $\Psi$

are estimated using an EM algorithm

8

# EM algorithm

## Data

- Observed data :  $\mathbf{x}_{1:n}$
- Missing (or hidden) data :  $\mathbf{z}_{1:n}$

## Principle

- Starting from  $\theta^0$
- At step  $q$ 
  - E(xpectation) step:  $Q(\theta, \theta^q) = E_{Z_{1:n}|\mathbf{x}_{1:n}}[\log P(\mathbf{x}_{1:n}, \mathbf{z}_{1:n}, \theta)]$
  - M(aximisation) step:  $\theta^{q+1} = \operatorname{argmax}_{\theta} Q(\theta, \theta^q)$

9

## EM for factor analysis

Let us assume that  $\boldsymbol{\mu} = \mathbf{0}$  (centering of the  $\mathbf{x}_i$ ), the complete log-likelihood is

$$\begin{aligned}\log p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \mathbf{W}, \boldsymbol{\Psi}) &= \sum_i \log \mathcal{N}_L(\mathbf{z}_i; \mathbf{0}, \mathbf{I}) + \log \mathcal{N}_D(\mathbf{x}_i; \mathbf{W}\mathbf{z}_i, \boldsymbol{\Psi}) \\ &= -\frac{n}{2} \log |\mathbf{I}_L| - \frac{n}{2} \text{Tr}(\hat{\boldsymbol{\Sigma}}_{zz}) \\ &\quad - \frac{n}{2} \log |\boldsymbol{\Psi}| - \frac{n}{2} \text{Tr}(\hat{\boldsymbol{\Sigma}}_{xx} \boldsymbol{\Psi}^{-1}) + \text{Cst}\end{aligned}$$

where

$$\hat{\boldsymbol{\Sigma}}_{xx} = \frac{1}{n} \sum_i (\mathbf{x}_i - \mathbf{W}\mathbf{z}_i)(\mathbf{x}_i - \mathbf{W}\mathbf{z}_i)^T$$

### Exercise

Demonstrate the above formula

10

## E step

The expectation of the complete log-likelihood requires

1.  $\mathbb{E}_{z|x}[\mathbf{z}_i] = \mathbf{S}\mathbf{W}^T\boldsymbol{\Psi}^{-1}(\mathbf{x}_i - \boldsymbol{\mu})$  where  $\mathbf{S} = (\mathbf{I}_L + \mathbf{W}^T\boldsymbol{\Psi}^{-1}\mathbf{W})^{-1}$
2.  $\mathbb{E}_{z|x}[\mathbf{z}_i\mathbf{z}_i^T] = \mathbb{E}_{z|x}[\mathbf{z}_i]\mathbb{E}_{z|x}[\mathbf{z}_i^T] + \mathbf{S}$

11

## M step

Reminders

$$\frac{\partial(b^T a)}{\partial a} = b$$

$$\frac{\partial(a^T A a)}{\partial a} = (A + A^T)a$$

$$\frac{\partial}{\partial A} \text{tr}(BA) = B^T$$

$$\frac{\partial}{\partial A} \log |A| = (A^{-1})^T$$

$$\text{tr}(ABC) = \text{tr}(CAB) = \text{tra}(BCA)$$

Thus if  $\mathbf{x}$  is a vector

$$\mathbf{x}^T A \mathbf{x} = \text{tr}(\mathbf{x}^T A \mathbf{x}) = \text{tr}(A \mathbf{x} \mathbf{x}^T)$$

12

## M step for $\Psi$

$$\mathbb{E}_{z|x} \left[ \frac{\partial L(\mathbf{W}, \Psi)}{\partial \Psi^{-1}} \right] = \mathbb{E}_{z|x} \left[ \frac{n}{2} (\Psi - \hat{\Sigma}_{xx}) \right] = \frac{n}{2} (\Psi - \mathbb{E}_{z|x} [\hat{\Sigma}_{xx}]) = 0$$

where

$$\begin{aligned} \mathbb{E}_{z|x} [\hat{\Sigma}_{xx}] &= \frac{1}{n} \left( \sum_i \mathbf{x}_i \mathbf{x}_i^T + W \left( \sum_i \mathbb{E}_{z|x} [z_i z_i^T] \right) W^T - 2W \sum_i \mathbf{E}_{z|x} [z_i] \mathbf{x}_i^T \right) \\ &= \frac{1}{n} \left( \sum_i \mathbf{x}_i \mathbf{x}_i^T + W \left( \sum_i \mathbb{E}_{z|x} [z_i z_i^T] \right) - 2W \sum_i \mathbf{E}_{z|x} [z_i] \mathbf{x}_i^T \right) \\ &= \frac{1}{n} \left( \sum_i \mathbf{x}_i \mathbf{x}_i^T - W \sum_i \mathbf{E}_{z|x} [z_i] \mathbf{x}_i^T \right) \end{aligned}$$

13

## M step for $\mathbf{W}$

$$\mathbb{E}_{z|x} \left[ \frac{\partial L(\mathbf{W}, \Psi)}{\partial \mathbf{W}} \right] = \mathbb{E}_{z|x} \left[ -\Psi^{-1} \sum_i \mathbf{x}_i \mathbf{z}_i^T + \Psi^{-1} \mathbf{W} \sum_i \mathbf{z}_i \mathbf{z}_i^T \right] = 0$$

14

## M Step summary

### Loading matrix

$$\mathbf{W}^{q+1} = \left( \sum_i (\mathbf{x}_i - \bar{\mathbf{x}}) \mathbb{E}_{z|x} [\mathbf{z}_i]^T \right) \left( \sum_i \mathbb{E}_{z|x} [\mathbf{z}_i \mathbf{z}_i^T] \right)^{-1}$$

### Noise covariance matrix

$$\Psi^{q+1} = \frac{1}{N} \text{diag} \left\{ \sum_i \mathbf{x}_i \mathbf{x}_i^T - \mathbf{W}^{q+1} \mathbb{E}_{z|x} [\mathbf{z}_i] \mathbf{x}_i^T \right\}$$

### Log-likelihood

The log-likelihood can be computed using the EM decomposition

$$\log P(\mathbf{X}; \Theta) = E_{Z_{1:n} | \mathbf{x}_{1:n}} [\log P(\mathbf{x}_{1:n}, \mathbf{z}_{1:n}; \theta)] - E_{Z_{1:n} | \mathbf{x}_{1:n}} [\log P(\mathbf{z}_{1:n} | \mathbf{x}_{1:n}; \theta)]$$

15



# Implementation of the algorithm

## Initialisation via a PCA

```
1 initialisation.FA<-function(X,L=1){
2   # Return W and Psi
3   d<-ncol(X)
4   Sigmaxx<-var(X)
5   W<-eigen(Sigmaxx)$vectors[,1:L]
6   if (L==1) W<-cbind(W)
7   Psi<-rep(1,d)
8   return(list(W=W,Psi=Psi))
9 }
```

## E step

```
1 FA.E.step<-function(X,W,Psi){
2   # X is assumed to be centered
3   # M contain the contionnal expectation of the latent factor
4   # S contains the covariance of the latent factor
5   L<-ncol(W)
6   S <- solve(diag(L) + t(W)%*%diag(1/Psi)%*%W)
7   M<- X%*%diag(1/Psi)%*%W%*%S
8   return(list(S=S,M=M))
9 }
```

## M Step

```
1 FA.M.step<-function(X,S,M,W,Psi){
2   n<-nrow(X)
3   Psi<-1/n*diag(t(X)%*%X -W%*%t(M)%*%X)
4   W<- (t(X)%*%M)%*%solve(n*S+t(M)%*%M)
5   return(list(Psi=Psi,W=W))
6 }
```

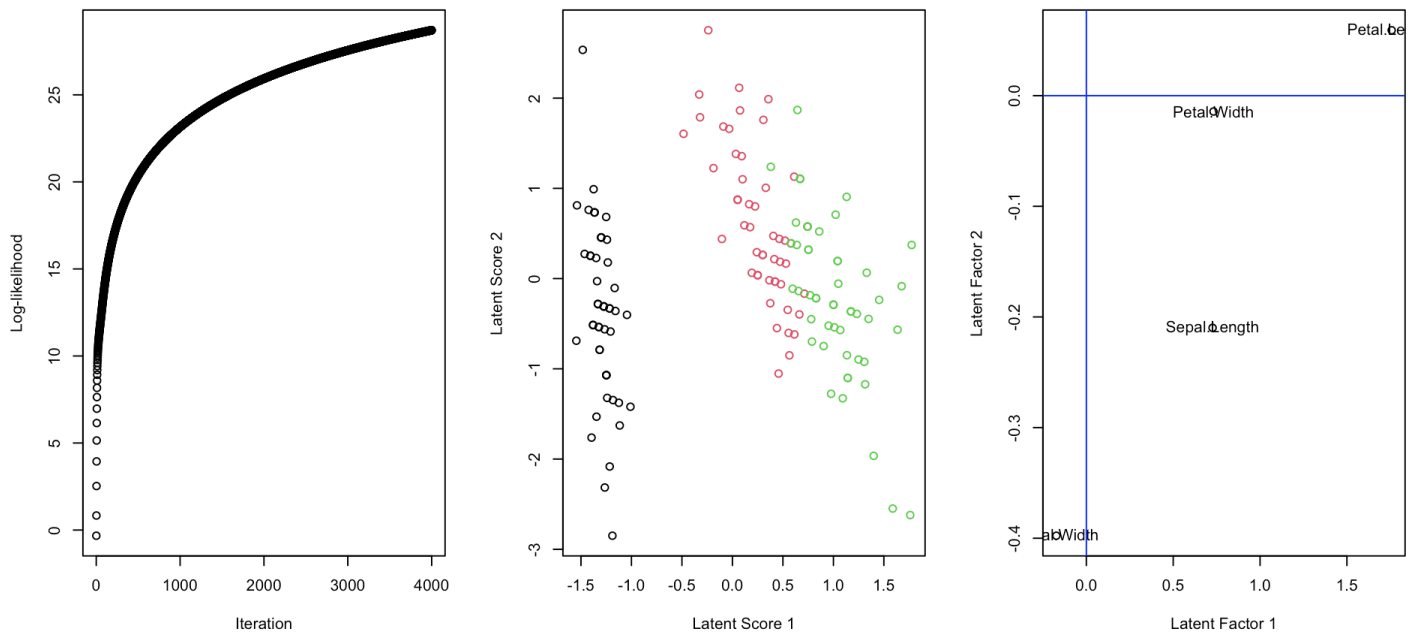
## Computation of the criterion

```
1 log.likelihood.FA<-function(X,S,M,Psi,W){
2   n<-nrow(X)
3   Sigmax<-(t(X)%*%X-W%*%t(M)%*%X)
4   return(-(sum(diag(S+t(M)%*%M/n))
5     +log(det(diag(Psi))))+
6     log(det(S))+
7     1/n*sum(diag(Sigmax%*%diag(1/Psi))))
8 }
```

## Putting it all together

```
1 FA.EM<-function(X,L=1,max.iter=50){
2   X<-scale(X, scale=FALSE);mu<-attr(X,"scaled:center")
3   log.likelihood<-NULL; init<-initialisation.FA(X,L)
4   W<-init$W; Psi<-init$Psi; criterion<- Inf; iteration<-1;
5   log.likelihood[iteration]<--Inf
6   while ((criterion>1e-6)&&(iteration<=max.iter)){
7     E.step<-FA.E.step(X,W,Psi); E.step$S->S; E.step$M->M
8     M.step<-FA.M.step(X,S,M,W,Psi); M.step$Psi->Psi; M.step$W->W
9     iteration<-iteration+1
10    log.likelihood[iteration]<-log.likelihood.FA(X,S,M,Psi,W)
11    criterion<-abs((log.likelihood[iteration] - log.likelihood[iteration-1])/max(log.likelihood[iteration], log.likelihood[iteration-1]))
12  }
13  return(list( W=data.frame(W), Psi=Psi,
14             M=data.frame(M), S=S,mu=mu,
```

## Example with the Iris



17

## Unidentifiability

- If we consider  $\mathbf{R}$  an orthogonal rotation matrix such that

$$\mathbf{R}\mathbf{R}^T = \mathbf{I}$$

It appears that  $\tilde{\mathbf{W}} = \mathbf{W}\mathbf{R}$  produces the same log-likelihood.

- $\mathbf{W}$  cannot be uniquely identified.

18

## Possible rotations

1. Forcing  $W$  to be orthogonal with columns ordered by decreasing variance
2. Forcing  $W$  to be lower triangular (problem of founder variables)
3. Choosing an informative rotation matrix. For example the varimax rotation.
4. ...

### Varimax

Varimax rotation maximizes the sum of the variance of the squared correlations between variables and factors

$$R_{\text{VARIMAX}} = \arg \max_R \left( \frac{1}{p} \sum_{j=1}^k \sum_{i=1}^p (WR)_{ij}^4 - \sum_{j=1}^k \left( \frac{1}{p} \sum_{i=1}^p (WR)_{ij}^2 \right)^2 \right)$$

This results in high factor loadings for a small number of variables and low factor loadings for the rest.

19

### Varimax

```
1 W.FA<-FA.result$W
2 W.Varimax<-varimax(as.matrix(FA.result$W))$loadings
3 print(W.Varimax)
```

Loadings:

	Latent Factor 1	Latent Factor 2
Sepal.Length	0.756	
Sepal.Width		-0.429
Petal.Length	1.683	0.509
Petal.Width	0.711	0.174

	Latent Factor 1	Latent Factor 2
SS loadings	3.916	0.473
Proportion Var	0.979	0.118
Cumulative Var	0.979	1.097

20

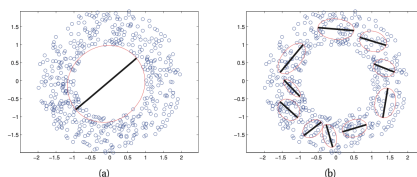
# Mixture of factor analysers

- Factor analyses is a way to estimate a variance matrix with few parameters
- This property can be used in the context of Gaussian mixture model assuming the following parameterization for component densities:

$$p(\mathbf{x}_i | \mathbf{z}_i, q_i = k) = \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k + \mathbf{W}_k \mathbf{z}_i, \boldsymbol{\Psi})$$

where  $k$  is the component number and  $\mathbf{W}_k$  is a loading matrix defining the relation between the observation  $\mathbf{x}_i$  and the latent vector  $\mathbf{z}_i$

- This approach is similar to the Banfield-Raftery idea of decomposing the component variance matrix  $\mathbf{k}$  in volume, form et direction.



Mixture of 1d PPCAs with 1 and 10 components (from Murphy Chapter 12)

## Relation to principal component analysis

### Assumption

If

- $\boldsymbol{\Psi} = \sigma^2 \mathbf{I}$
- $\mathbf{W}$  is orthogonal

and

- $\sigma^2 \rightarrow 0$

Then

Tipping, M. and C. Bishop (1999, Probabilistic principal component analysis. J. of Royal Stat. Soc. Series B 21(3), 611–622) showed that FA is equivalent to PCA

### Criterion

$$J(\mathbf{W}, \mathbf{Z}) = \|\mathbf{X} - \mathbf{Z}\mathbf{W}^T\|_F^2$$

where  $\mathbf{W}^T\mathbf{W} = \mathbf{I}$

22

## A Constrained EM Algorithm for PCA (from Ahn, J.-H. and J.-H. Oh, 2003)

```

1 upper<-function(A){A[lower.tri(A,diag=FALSE)]<-0;return(A)}
2 lower<-function(A){A[upper.tri(A,diag=FALSE)]<-0;return(A)}
3 PCA.EM<-function(X,q=2){
4   p<-ncol(X); n<-nrow(X)
5   W<-diag(p)[,1:q]; M<-X%%W # Initialisation
6   Jold<-0; J<-1; iteration<-0; Error<-NULL
7   while ((abs(J - Jold)>1e-3)){
8     Jold<-sum((X-M%%t(W))^2)
9     S <- solve(upper(t(W)%%W)); M<- X%%W%%S # E-step
10    W<- (t(X)%%M)%%solve(lower(n*S+t(M)%%M)) # M-step
11    W<-apply(W,2,function(x){x/sqrt(sum(x^2))})#orthogonalisation
12    J<-sum((X-M%%t(W))^2); Error[iteration<-iteration+1]<-J
13  }
14  return(list(W=data.frame(W),M=data.frame(M),Error=Error))}

```

23

# Independent Component Analysis

24

## Independent Component analysis (Wikipedia)

Independent component analysis attempts to decompose a multivariate signal into independent non-Gaussian signals.

### *Cocktail party problem*

Underlying speech signals are separated from a sample data consisting of people talking simultaneously in a room.

That the ICA separation of mixed signals gives very good results is based on two assumptions

### *Two assumptions:*

- The source signals are independent of each other.
- The values in each source signal have non-Gaussian distributions.

25

## ICA Historical context (French Wikipedia)

### *Blind source separation*

The first formulation was carried out in 1984 by Jeanny Hérault and Bernard Ans, two researchers in neuroscience and signal processing, to model in the form of a neuromimetic network self-adaptive encoding and decoding of movement in humans.

### *France and Finland*

- The French signal processing community adopted a statistical formalism
- While Finnish researchers aimed to extend principal component analysis by means of a connectionist formalism (1985)

### *Formalisation*

26

## ICA Model

Let  $\mathbf{x}_t \in \mathbb{R}^D$  be the observed signal at the sensors at "time"  $t$ , and  $\mathbf{z}_t \in \mathbb{R}^L$  be the vector of source signals:

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

- $\mathbf{W}$  is an  $D \times L$  matrix,
- $\boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Psi})$ .

The model is identical to factor analysis (or PCA if there is no noise, except we don't in general require orthogonality of  $\mathbf{W}$ ).

However, we will use a **different prior** for  $p(\mathbf{z}_t)$ .

In FA, we assume each source is independent, and has a Gaussian distribution.

### *Relax this Gaussian assumption on latent variables (sources)*

$$p(\mathbf{z}_t) = \prod_j p_j(z_{tj}).$$

### Additional assumptions

- Without loss of generality the variance of the source distributions is constrained to unity
- $\mathbf{W}$  is assumed square and hence invertible.

27

## Maximum likelihood estimation of ICA

If the data is centered and whitened, we have

$$\mathbb{E}[\mathbf{x}\mathbf{x}^t] = \mathbf{I} = \mathbf{W}\mathbb{E}[\mathbf{z}\mathbf{z}^T]\mathbf{W}^T = \mathbf{W}\mathbf{W}^T$$

also have

Hence we see that  $\mathbf{W}$  must be orthogonal. This reduces the number of parameters we have to estimate from  $D^2$  to  $D(D - 1)/2$ .

### Recognition weights/Generative weights

Let  $\mathbf{V} = \mathbf{W}^{-1}$  these are often called the recognition weights, as opposed to  $\mathbf{W}$ , which are the generative weights

### Log-likelihood

Since  $\mathbf{x} = \mathbf{W}\mathbf{z}$ , we have

$$p_x(\mathbf{W}\mathbf{z}_t) = p_z(\mathbf{z}_t)|\det(\mathbf{W}^{-1})| = p_z(\mathbf{V}\mathbf{x}_t)|\det(\mathbf{V})|$$



Hence assuming T iid samples:

$$L(\mathbf{V}) = \frac{1}{T} \log p(\mathcal{D}|\mathbf{V}) = \underbrace{\log |\det(\mathbf{V})|}_0 + \frac{1}{T} \sum_{jt} \log p_j(v_j^T \mathbf{x}_t)$$

since  $\mathbf{V}$  is orthogonal.

## Estimation via Gradient Ascent

Let us define  $h_j = \mathbf{v}_j^T \mathbf{x}$ ,

$$g_j(h_j) = \frac{\partial \log p_j(h_j)}{\partial h_j},$$

$$\frac{\partial L(\mathbf{V})}{\partial V_{ij}} = W_{ji} + x_i g_j(h_j).$$

where the datapoint  $\mathbf{x}^T = (x_i)_{i=1 \dots D}$ .

Repeat for each datapoint  $\mathbf{x}$ :

1. Put  $\mathbf{x}$  through a linear mapping:

$$\mathbf{h} = \mathbf{V}\mathbf{x}$$

2. Put  $\mathbf{h}$  through a nonlinear map:

$$g_j = g_j(h_j)$$

, where a popular choice is  $g() = -\tanh()$ .

3. Adjust the weights in accordance with

$$\nabla \mathbf{V} \propto [\mathbf{V}^T]^{-1} + \mathbf{x}\mathbf{g}^T.$$

- matrix inversion results in a slow algorithm

## Fast ICA (for one latent factor $\mathbf{v}$ )

Let consider  $G(z) = -\log p(z)$ ,  $g(z) = G'(z)$ .

- $L(\mathbf{v}) = \mathbb{E}[G(\mathbf{v}^T \mathbf{x})] + \lambda(1 - \mathbf{v}^T \mathbf{v})$ , the theoretical objective function (**to be minimized**)
- $\nabla L(\mathbf{v}) = \mathbb{E}[\mathbf{x}g(\mathbf{v}^T \mathbf{x})] - 2\lambda\mathbf{v}$ , the gradient
- $H(\mathbf{v}) = \mathbb{E}[\mathbf{x}\mathbf{x}^T g'(\mathbf{v}^T \mathbf{x})] - 2\lambda I$ , the hessian matrix

Let us make the approximation

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T g'(\mathbf{v}^T \mathbf{x})] = \underbrace{\mathbb{E}[\mathbf{x}\mathbf{x}^T]}_I \mathbb{E}[g'(\mathbf{v}^T \mathbf{x})] = \mathbb{E}[g'(\mathbf{v}^T \mathbf{x})]$$

This makes the Hessian very easy to invert, giving rise to the following Newton update:

$$\mathbf{v}^* \triangleq \mathbf{v} - H(\mathbf{v})^{-1} \nabla L(\mathbf{v}) = \mathbf{v} - \frac{\mathbb{E}[\mathbf{x}g(\mathbf{v}^T \mathbf{x})] - 2\lambda\mathbf{v}}{\mathbb{E}[g'(\mathbf{v}^T \mathbf{x})] - 2\lambda}$$

Which can be expressed as

$$\mathbf{v} := E[\mathbf{x}g(\mathbf{v}^T \mathbf{x})] - E[g'(\mathbf{v}^T \mathbf{x})]\mathbf{v}$$

(In practice, the expectations can be replaced by Monte Carlo estimates from the training set, which gives an efficient online learning algorithm.)

After performing this update, one should project back onto the constraint surface using

$$\mathbf{v} := \frac{\mathbf{v}}{\|\mathbf{v}\|}$$

30

## Fast ICA for $L > 1$

Centering and Withining is assumed

For  $L > 1$  the process is iterated for all  $\mathbf{v}_j$  with orthogonalisation

### Fast ICA for $L > 1$

- Intitilisation of  $\mathbf{V}$
- **for**  $j$  in 1 to  $L$ :
  - **while**  $\mathbf{v}_j$  changes
    - ↪ Newton update:  $\mathbf{v}_j := E[\mathbf{x}g(\mathbf{v}_j^T \mathbf{x})] - E[g'(\mathbf{v}_j^T \mathbf{x})]\mathbf{v}_j$
    - ↪ Gram-Schmidt orthogonalisation:  $\mathbf{v}_j := \mathbf{v}_j - \sum_{k=1}^{j-1} (\mathbf{v}_j^T \mathbf{v}_k) \mathbf{v}_k$
    - ↪ Normalisation:  $\mathbf{v}_j := \frac{\mathbf{v}_j}{\|\mathbf{v}_j\|}$
- Output:  $\mathbf{Z} = \mathbf{XV}$

31

## Non-Gaussianity

For non-Gaussianity, FastICA relies on a nonquadratic nonlinear function  $f(u)$ , its first derivative  $g(u)$ , and its second derivative  $g'(u)$ .

### *Classical cost*

- $G(u) = \log \cosh(u)$  which gives  $g(u) = \tanh(u)$
- $G(u) = -\exp(-u^2/2)$  which gives  $g(u) = u \exp(-u^2/2)$
- $G(u) = u^4/4$  which gives  $g(u) = u^3$

Neural networks and unsupervised  
learning

# Modeling of a neuron

The first modeling of the neuron was suggested in the 1940s by Mac Culloch and Pitts. It was a unit which according to several signals transmitted a binary response.

In general, a formal neuron has

- dendrites which receive the input signal and
- an axon which transmits the output signal

## *The input signal*

$\boldsymbol{x}$  is a vector belonging most often to  $\mathbb{R}^d$  or  $\{0, 1\}^d$ .

## *Dendrites*

are characterized by a weight vector  $\boldsymbol{w}$

## *The output*

is a function of  $\boldsymbol{x}$  and  $\boldsymbol{w}$ , which is the composition of an input function,  $h(\boldsymbol{x}, \boldsymbol{w})$  and an output (or activation) function,  $f(h)$

## The neuron as a function

Most of the time the input function is a simple dot product:

$$h(\mathbf{x}, \mathbf{w}) = \mathbf{x}^T \mathbf{w}$$

The activation functions are diverse but belong to large families (radial bases, sigmoid functions ...).

A typical activation function is for example:

$$f(\mathbf{x}, \mathbf{w}) = \eta \cdot \frac{\exp \{(\mathbf{x}, \mathbf{w})\} - 1}{\exp(\mathbf{x}, \mathbf{w}) + 1}.$$

The functions which have this appearance are said to be sigmoid

35

## Neural networks

- Neurons can be connected to each other and then form a network.
- Learning consists of adjusting the free network parameters according to the desired goal, that is, to calculate the values of the weight vectors as a function of the inputs.

### *Two layers networks*

- with one input layer transmitting the input vector to the second layer neurons
- with one output layer compressing the information of the first layer with linear activation function
- allows to rewrite k-means and PCA with online learning (gradient descent)

36

## Hebbian Learning (Hebb 1949)

### Principle

An increase of synaptic strength between an input and an output neuron may be related to the firing rates of the input and output [Hebb, 1949].

### Practical implementation

As a result, synaptic strengths will increase fastest between pairs of neurons whose responses are correlated, and the resulting increase in synaptic strength will lead to a further increase in the correlation.

$$\Delta \mathbf{w} = \eta y(\mathbf{x}) \mathbf{x},$$

or in scalar form with implicit  $n$ -dependence,

$$w_i(n+1) = w_i(n) + \eta y(\mathbf{x}) x_i$$

Increasing the correlation in this manner may lead to a useful pattern of synaptic strengths over a population of neurons.

37

## Stochastic Gradient Descent (from Wikipedia)

Statistical estimation and machine learning consider the problem of minimizing an objective function that has the form of a sum:

$$Q(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n Q_i(\mathbf{w}),$$

where the parameter  $\mathbf{w}$  that minimizes  $Q(\mathbf{w})$  is to be estimated.

Each summand function  $Q_i$  is typically associated with the  $i$  – *th* observation in the data set (used for training).

Sum-minimization problems arise:

- in least squares and in maximum-likelihood estimation,
- empirical risk minimization.

When used to minimize the above function, a standard (or “batch”) gradient descent method would perform the following iterations:

$$\mathbf{w} := \mathbf{w} - \eta \nabla Q(\mathbf{w}) = \mathbf{w} - \frac{\eta}{n} \sum_{i=1}^n \nabla Q_i(\mathbf{w}),$$

where  $\eta$  is a step size (sometimes called the learning rate in machine learning).

38

## Iterative method

Fluctuations in the total objective function as gradient steps with respect to mini-batches are taken.

In stochastic (or “on-line”) gradient descent, the true gradient is approximated by a gradient at a single example:

$$\mathbf{w} := \mathbf{w} - \eta \nabla Q_i(\mathbf{w}).$$

- As the algorithm sweeps through the training set, it performs the above update for each training example.
- Several passes can be made over the training set until the algorithm converges.
- If this is done, the data can be shuffled for each pass to prevent cycles.
- Typical implementations may use **an adaptive learning rate** so that the algorithm converges.

39



## Stochastic Gradient in pseudocode

1. Choose an initial vector of parameters  $\mathbf{w}$  and learning rate  $\eta$
2. Repeat until an approximate minimum is obtained:
  - a. Randomly shuffle examples in the training set.
  - b. For  $i = 1, 2, \dots, n$  do:  $\mathbf{w} := \mathbf{w} - \eta \nabla Q_i(\mathbf{w})$ .

40

## Adaptative learning rate

A distinction exists between constant gain algorithms,

$$\eta_n \geq 0, \quad \lim_{n \rightarrow \infty} \eta_n = \eta > 0$$

and decreasing gain algorithms,

$$\sum_{n=0}^{\infty} \eta_n = \infty, \quad \sum_{n=0}^{\infty} \eta_n^2 < \infty.$$

The first are dedicated to the estimation of parameters changing slowly over time and the second to the estimation of stable parameters.

41

## K-means and Winner take all

is a computational principle applied in computational models of neural networks by which neurons in a layer compete with each other for activation.

The simplest form of competitive learning modifies only the weight vector of “the best” neuron at every stage of learning.

In fact, with each presentation of a input (a vector of the training set), two steps are performed:

1. choose the best neuron, i.e. the one that shows the most important output
2. modify the weight vector of this neuron.

When the activation function is increasing (which is not true for the functions with radial basis), the winning neuron is the one that produces the greatest value of function entry.

If we consider a dot product as an input function, the weight vector of the winner,  $i^*$ , checks:

$$\forall i, (\mathbf{w}_{i^*} \cdot \mathbf{x}) \geq (\mathbf{w}_i \cdot \mathbf{x}).$$

And if *the weight vectors are normalized*, the winner is the neuron that has the weight vector, closest to the input  $\mathbf{x}$ , in the sense of the Euclidean distance.

The coordinates of the winner’s weight vector are updated using a rule of the type following :

$$\mathbf{w}_{i^*}(t + 1) = \mathbf{w}_{i^*}(t) + \eta(t) \cdot (\mathbf{x} - \mathbf{w}_{i^*}(t)), \eta(t) \leq 1$$

where  $\eta(t)$  is the training step at iteration  $t$ .

# Stochastic Gradient for Kmeans (online Kmeans)

The criterion to be optimized can be written as

$$Q(\mathbf{w}_1, \dots, \mathbf{w}_K) = \frac{1}{2n} \sum_i \sum_k I_{(k=\arg \min_{\ell} \|\mathbf{x}_i - \mathbf{w}_{\ell}\|^2)} \|\mathbf{x}_i - \mathbf{w}_k\|^2$$

1. Choose an initial vector of parameters  $\mathbf{w}$  and learning rate  $\eta$
2. Repeat until an approximate minimum is obtained:
  - a. Randomly shuffle examples in the training set.
  - b. For  $i = 1, 2, \dots, n$  do:
    1. For  $k = 1, 2, \dots, K$ 
      - $\nabla Q_i(\mathbf{w}_k) = -\frac{1}{n} I_{(k=\arg \min_{\ell} \|\mathbf{x}_i - \mathbf{w}_{\ell}\|^2)} (\mathbf{x}_i - \mathbf{w}_k)$
      - $\mathbf{w}_k := \mathbf{w}_k - \eta \nabla Q_i(\mathbf{w}_k)$ .

43

## Kmeans implementation

```
1 kmeans.winner.take.all<-function(X,K=2,max.iteration=2000){
2   p<-ncol(X); n<-nrow(X);shuffling<-sample(1:n,n)
3   X<-X[shuffling,]; W<-X[sample(1:n,K),]
4   Q<-rep(0,max.iteration); cluster<-rep(0,n)
5   distances<-rep(sum(diag(var(X)))*(n-1)/n,n)
6   for (i in 1:max.iteration){
7     x<-cbind(X[(i-1)%n + 1,])
8     distances.x.to.W<-sum(x^2)-2*as.matrix(x)%*%t(W)+ colSums(t(W^2))
9     winner.index<-which.min(distances.x.to.W)
10    W[winner.index,]<-W[winner.index,] + 1/i*(x-W[winner.index,])
11    cluster[(i-1)%n + 1]<-winner.index
12    distances[(i-1)%n + 1]<-distances.x.to.W[winner.index]
13    Q[i]<-mean(distances) }
14   return(list(W=W,Q=Q,cluster=cluster[order(shuffling)]))
15 }
```

44

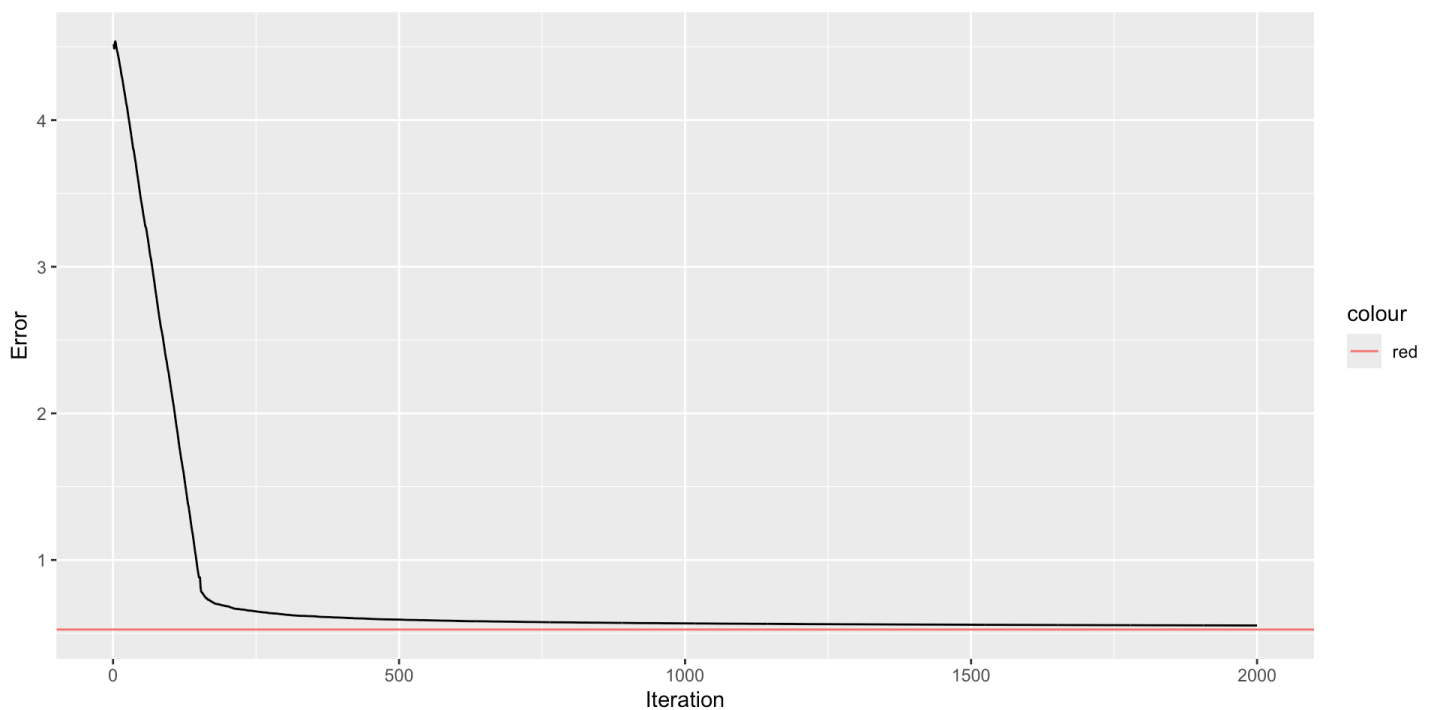
## One line kmeans example with Fisher iris

```
1 data(iris)
2 X<-iris[,1:4]
3 set.seed(1)
4 kmeans.winner.take.all(X,3)->res
5 table(res$cluster,iris$Species)
```

	setosa	versicolor	virginica
1	50	0	0
2	0	45	11
3	0	5	39

45

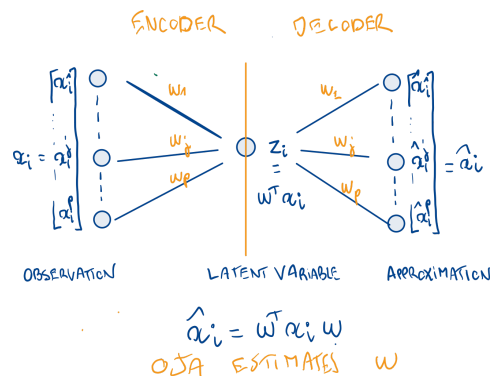
## One line kmeans example with Fisher iris



46

## PCA and Oja's rule

Consider a linear neuron with output  $z = \mathbf{w}^T \mathbf{x}$  that returns a linear combination of its inputs  $\mathbf{x}$  using presynaptic weights  $\mathbf{w}$ .



PCA according OJA

Oja's rule defines the change in presynaptic weights  $\mathbf{w}$  given the output response  $y$  of a neuron to its inputs  $\mathbf{x}$  to be

$$\mathbf{w} := \mathbf{w} - \eta z (\mathbf{x} - z \mathbf{w})$$

47

## Stochastic Gradient PCA and Oja's rule

The criterion to be optimized can be written as

$$Q(\mathbf{w}) = \frac{1}{2n} \sum_i \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2 = \frac{1}{2n} \sum_i \|\mathbf{x}_i - y_i \mathbf{w}\|^2 = \frac{1}{2n} \sum_i \|\mathbf{x}_i - \mathbf{w}^T \mathbf{x}_i \mathbf{w}\|^2$$

where  $\|\mathbf{w}\|^2 = 1$ .

1. Choose an initial vector of parameters  $\mathbf{w}$  and learning rate  $\eta$
2. Repeat until an approximate minimum is obtained:
  - a. Randomly shuffle examples in the training set.
  - b. For  $i = 1, 2, \dots, n$  do:
    - $\nabla Q_i(\mathbf{w}) = y_i (\mathbf{x}_i - y_i \mathbf{w})$
    - $\mathbf{w} := \mathbf{w} - \eta \nabla Q_i(\mathbf{w})$ .

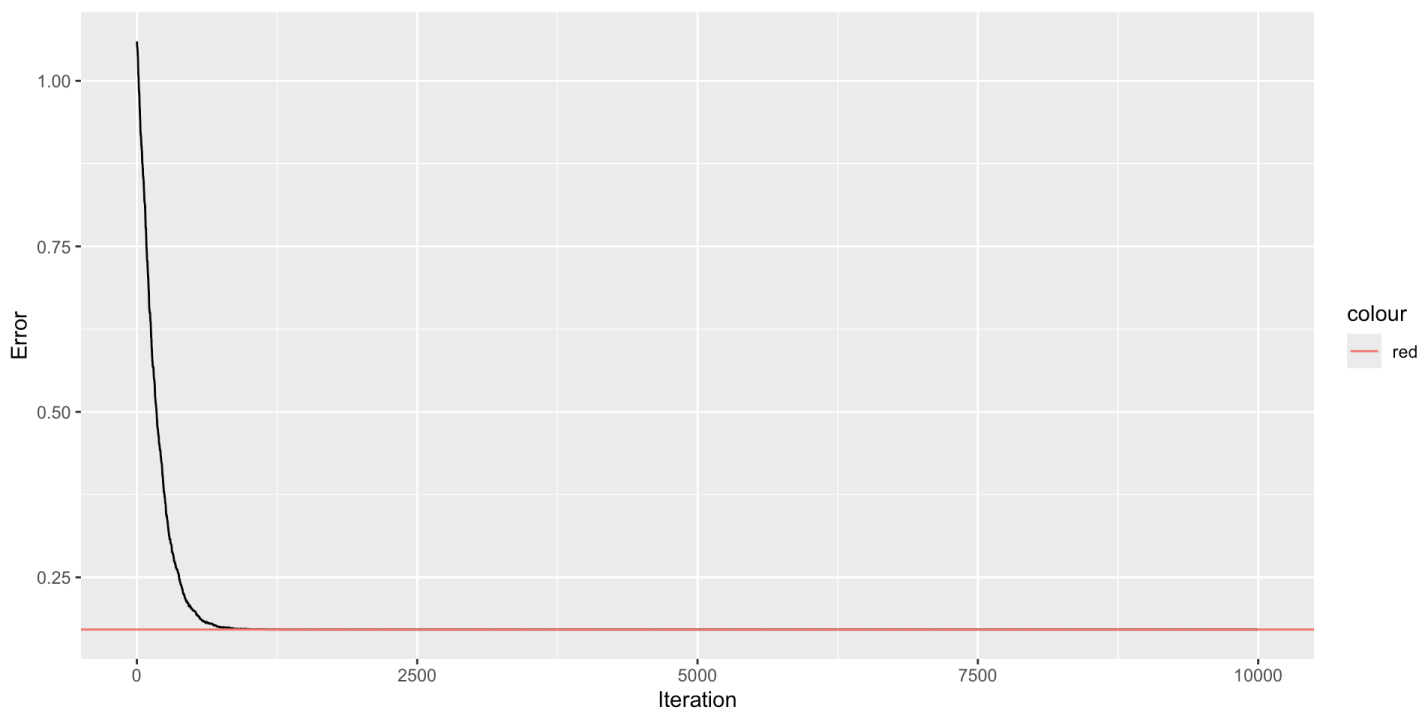
48

# Oja's rule implementation

```
1 Oja.rule<-function(X,max.iteration=10000,eta=0.001){
2   p<-ncol(X); n<-nrow(X)
3   w<-rbind(rep(1,p)); w<-w/(sqrt(sum(w^2)))
4   Q<-rep(0,max.iteration)
5   for (i in 1:max.iteration){
6     Q[i]<- 1/(2*n) *sum((X - X%%t(w)%*%w)^2)
7     x<-X[(i-1)%n + 1,]
8     y<-sum(w*x)
9     w<-w + eta*y*(x-w*y)}
10  return(list(w=w,Q=Q))
11 }
```

49

# Oja's rule example with Fisher iris



50

# Variational auto-encoder

51

## Variational autoencoder ideas

### *The original papers*

- Rezende, D. J., Mohamed, S., & Wierstra, D. (2014, June). Stochastic backpropagation and approximate inference in deep generative models. In International conference on machine learning (pp. 1278-1286). PMLR.
- Kingma, D. P., & Welling, M. (2013). Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114.

### *Main reference*

- Diederik P. Kingma and Max Welling (2019), "An Introduction to Variational Autoencoders", Foundations and Trends R in Machine Learning:

### *What is does*

- generate realistic samples of data,
- allow for accurate imputations of missing data,
- high-dimensional data visualisation
- Clustering

## How it works

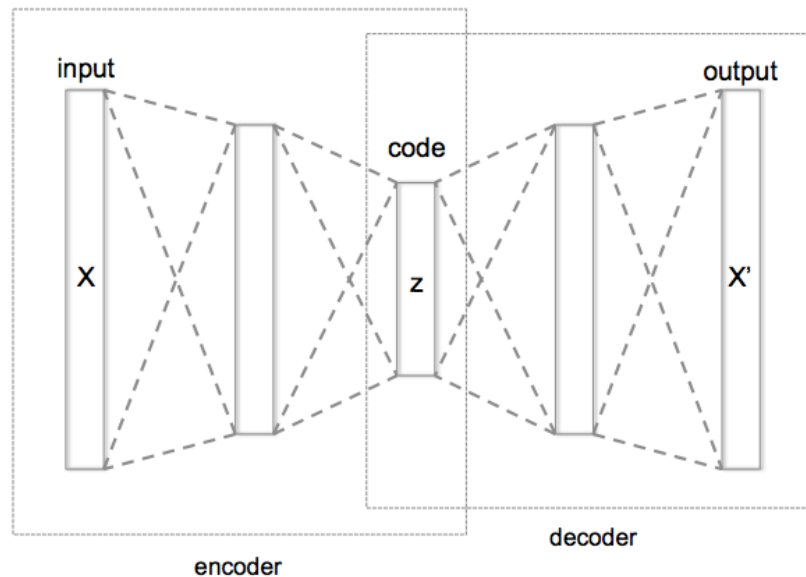
**Latent variables** models which marry ideas from

- approximate Bayesian inference
  - ELBO (Evidence Lower Bound)
  - Reparametrization
- deep neural networks
  - Stochastic Gradient Descent
  - Retropropagation of the Gradient

to represent an approximate posterior distribution through variational lower bound optimization

52

## Auto-encoder Structure





# What is a VAE ?

## Coupling of 2 parametric models

VAE is a latent variable vector  $\mathbf{z}$  and an observation  $\mathbf{x} \in \mathbb{R}^D$

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

- encoder  $g$  (recognition model):  $p(\mathbf{z}|\mathbf{x})$ , which is approximated by  $q_{\Phi}(\mathbf{z}|\mathbf{x})$
- decoder  $h \approx g^{-1}$  (generative model):  $p_{\Theta}(\mathbf{x}|\mathbf{z})$
- encoder and decoder could be neural networks

## Optimization of ELBO via Stochastic Gradient Ascent

- The VAE ELBO approximates the likelihood of a **latent variable** model
- The Gradient computation uses the re-parametrization trick
- Each step of the gradient ascent augment the ELBO as an **EM iteration**

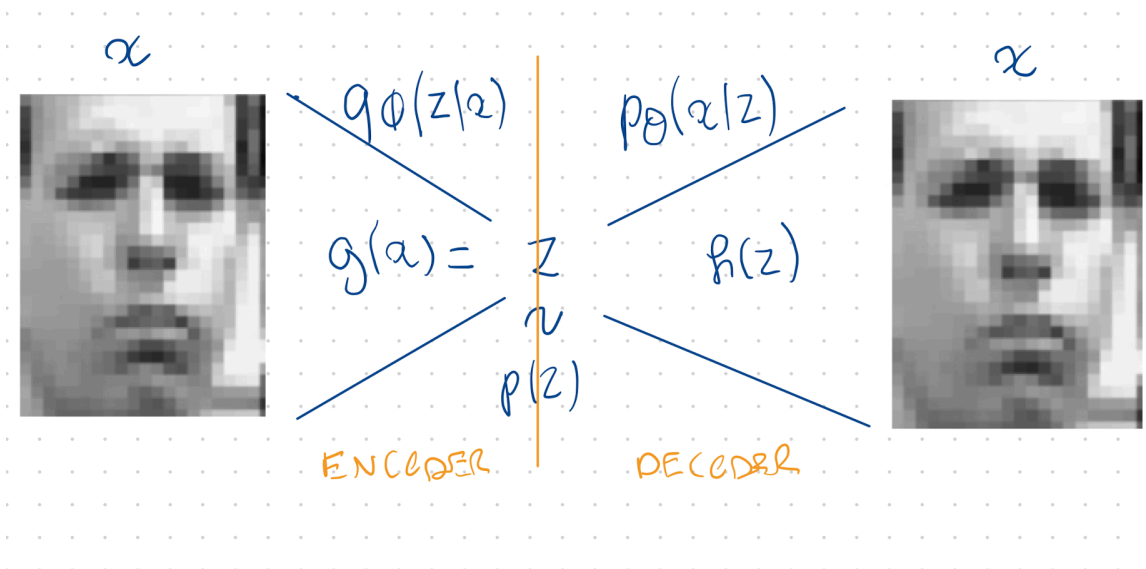
54

# Example of use

## Data

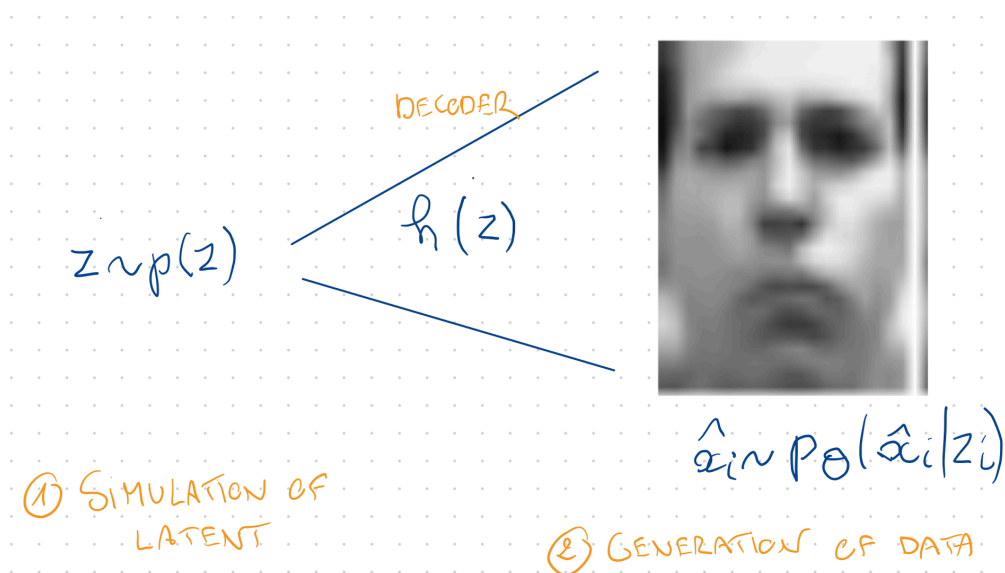
## Example of use

### Learning from data



# Example of use

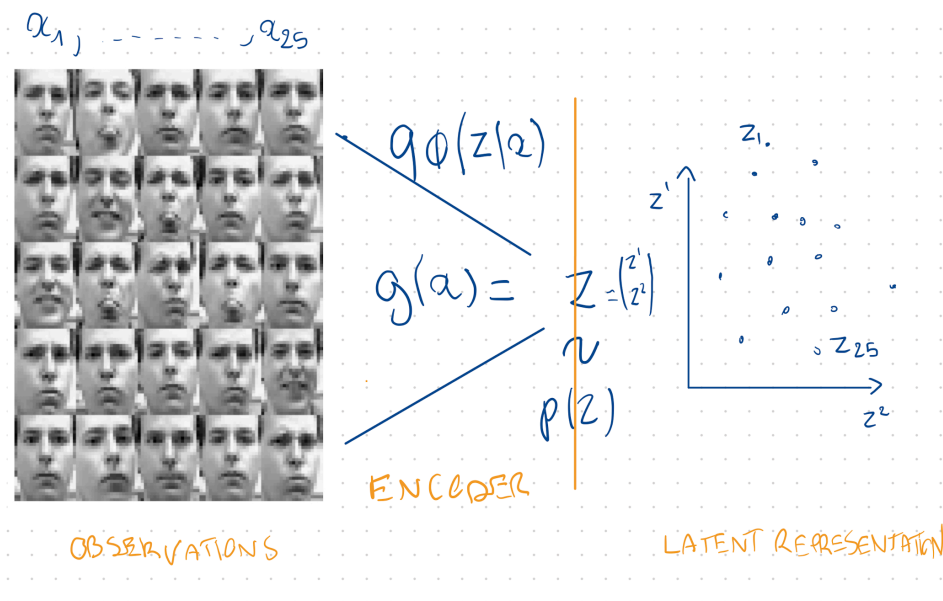
## New Data generation from latent simulation



Frey image generation

# Example of use

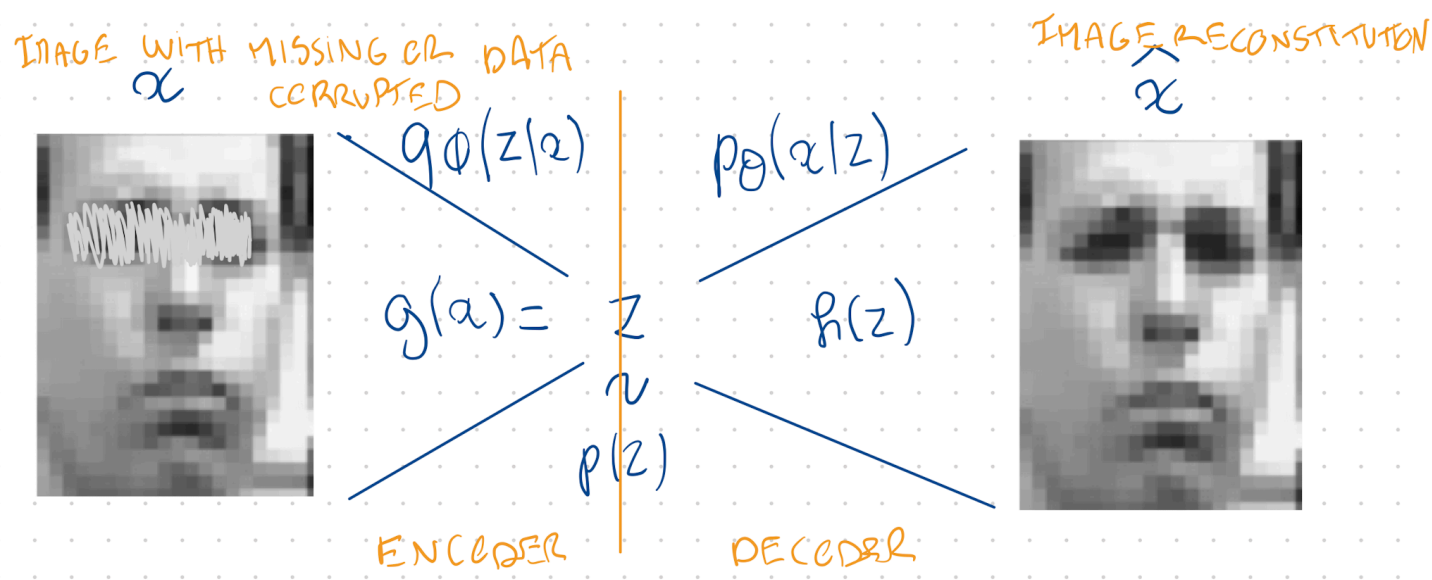
## Data representation in latent space



Frey image representation

## Example of use

### Missing data imputation after learning



Frey image pixel imputation

## Ingredient: Parameterization of conditional distributions with Neural Networks

## Modeling joint distribution

- $\mathbf{x}$ : Observed random variables
- $p^*(\mathbf{x})$ : underlying unknown distribution
- $p_\theta(\mathbf{x})$ : model distribution
- Goal:  $p_\theta(\mathbf{x}) \approx p^*(\mathbf{x})$

We wish flexible  $p_\theta(\mathbf{x})$

61

## Modeling Conditional distribution

*Classification and regression*

$$p_\theta(y|\mathbf{x}) \approx p^*(y|\mathbf{x})$$

62

# Parameterization of conditional distributions with Neural Networks

## Classification

$$\theta = \text{NeuralNet}(\mathbf{x})$$

$$p_{\theta}(y|\mathbf{x}) = \text{Categorical}(y, \theta)$$

## Ingredient: Stochastic Gradient

## What differences between Oja's rule and VAE

### What is a VAE ?

Coupling of 2 **parametric models**

- decoder (generative model):  $\mathbf{x} = h_{\Theta}(\mathbf{z})$  where  $p_{\Theta}(\mathbf{x}|\mathbf{z})$
- encoder (recognition model):  $\mathbf{z} = g_{\Phi}(\mathbf{x})$  where  $p(\mathbf{z}|\mathbf{x})$  is approximated by  $q_{\Phi}(\mathbf{z}|\mathbf{x})$

### In Oja's rule

- There are No probabilistic models
  - $z = g(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$  and
  - $\hat{\mathbf{x}} = h(z) = z\mathbf{w}$ ,
- Encoder and decoder share the same parameters  $\mathbf{w} = \Theta = \Phi$

65

## Factor analysis generalizes Oja and is closer to a VAE

Factor analysis considers the observation  $\mathbf{x} \in \mathbb{R}^D$

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

where

- the noise  $\boldsymbol{\epsilon} \sim \mathcal{N}_D(\mathbf{0}, \boldsymbol{\Psi})$
- the hidden (latent) vector  $\mathbf{z} \sim \mathcal{N}_L(\mathbf{0}, \mathbf{I}_L)$

$$p(\mathbf{x}|\mathbf{z}, \theta) = \mathcal{N}(\mathbf{x}; \mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \boldsymbol{\Psi})$$

the mean is a linear function of the (hidden) inputs

- $\mathbf{W}$  is a  $D \times L$  matrix, known as the factor loading matrix,
- $\boldsymbol{\Psi}$  is a  $D \times D$  covariance matrix that we take to be diagonal

The special case in which  $\boldsymbol{\Psi} = \sigma^2 \mathbf{I}$  is called probabilistic principal components analysis or PPCA.

66

### Missing data

In a missing data framework the log-likelihood of the parameters is advantageously expressed as

$$\log P_{\Theta}(\mathbf{x}) = \mathbb{E}_{\mathbf{z}|\mathbf{x}} \left[ \log \frac{P(\mathbf{x}, \mathbf{z})}{P(\mathbf{z}|\mathbf{x})} \right]$$

### Approximation

When the distribution of  $\mathbf{z}|\mathbf{x}$  is intractable, an **approximation**  $q_{\Phi}(\mathbf{z}|\mathbf{x})$  is used

$q_{\Phi}(\mathbf{z}|\mathbf{x})$  is the inverse function for  $p_{\Theta}(\mathbf{x}|\mathbf{z})$  in a **Bayes sense**



$$\begin{aligned}
\log P_{\Theta}(\mathbf{x}) &= \mathbb{E}_{q_{\Phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{P_{\Theta}(\mathbf{x}, \mathbf{z})}{P_{\Theta}(\mathbf{z}|\mathbf{x})} \right] \\
&= \mathbb{E}_{q_{\Phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{P_{\Theta}(\mathbf{x}, \mathbf{z}) q_{\Phi}(\mathbf{z}|\mathbf{x})}{P_{\Theta}(\mathbf{z}|\mathbf{x}) q_{\Phi}(\mathbf{z}|\mathbf{x})} \right] \\
&= \underbrace{\mathbb{E}_{q_{\Phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{P_{\Theta}(\mathbf{x}, \mathbf{z})}{q_{\Phi}(\mathbf{z}|\mathbf{x})} \right]}_{ELBO} + \underbrace{\mathbb{E}_{q_{\Phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{q_{\Phi}(\mathbf{z}|\mathbf{x})}{p_{\Theta}(\mathbf{z}|\mathbf{x})} \right]}_{D_{KL}(q_{\Phi}(\mathbf{z}|\mathbf{x})|p_{\Theta}(\mathbf{z}|\mathbf{x}))}
\end{aligned}$$

where  $D_{KL}(q_{\Phi}(\mathbf{z}|\mathbf{x})|p_{\Theta}(\mathbf{z}|\mathbf{x})) = -\mathbb{E}_q[\log \frac{p}{q}] \geq -\log \mathbb{E}_q[\frac{p}{q}] \geq 0$  from Jensen

## Two for one

- VAE finds parameters which approximately maximize the marginal likelihood  $P_{\Theta}(\mathbf{x})$  (good generative function)
- VAE finds the approximation of the recognition model which minimizes the KL divergence

## Alternative formulation of ELBO

ELBO can be rewritten as

$$ELBO = E_{q_{\Phi}(z|\mathbf{x})} [\log P_{\Theta}(\mathbf{x}|z)] - D_{KL}(q_{\Phi}(z|\mathbf{x})|p_{\Theta}(z))$$

For a suited choice of  $p(z)$  and  $q(z|\mathbf{x})$ ,  $D_{KL}(q_{\Phi}(z|\mathbf{x})|p_{\Theta}(z))$  can be calculated in closed form.

### Exercices

1. Show that the ELBO can be rewritten as above
2. Compute the KL divergence between two multivariate Gaussians

71

## The ELBO is maximized by Stochastic Gradient

Let  $\mathbf{X}$  be a i.i.d sample of random vector  $\mathbf{x}$

$$ELBO(\mathbf{X}) = \sum_{\mathbf{x} \in \mathbf{X}} \mathcal{L}(\Theta, \Phi; \mathbf{x})$$

### Gradient

The Gradient can be separated into 2 parts

1.  $\nabla_{\Theta} \mathcal{L}(\Theta, \Phi; \mathbf{x})$
2.  $\nabla_{\Phi} \mathcal{L}(\Theta, \Phi; \mathbf{x})$

72

## Decoder Gradient: $\nabla_{\Theta} \mathcal{L}(\Theta, \Phi; \mathbf{x})$

Given some usually verified conditions and a Monte Carlo Approximation

$$\begin{aligned}\nabla_{\Theta} \mathcal{L}(\Theta, \Phi; \mathbf{x}) &= \nabla_{\Theta} \mathbb{E}_{q_{\Phi}(z|\mathbf{x})} \left[ \log \frac{P_{\Theta}(\mathbf{x}, z)}{q_{\Phi}(z|\mathbf{x})} \right] \\ &= \nabla_{\Theta} \mathbb{E}_{q_{\Phi}(z|\mathbf{x})} [\log P_{\Theta}(\mathbf{x}, z)] \\ &= \mathbb{E}_{q_{\Phi}(z|\mathbf{x})} [\nabla_{\Theta} \log P_{\Theta}(\mathbf{x}, z)] \\ &\approx \nabla_{\Theta} \log P_{\Theta}(\mathbf{x}, z)\end{aligned}$$

73

## Encoder Gradient: $\nabla_{\Phi} \mathcal{L}(\Theta, \Phi; \mathbf{x})$

Encoder Gradient is more difficult to compute since in general

$$\begin{aligned}\nabla_{\Phi} \mathcal{L}(\Theta, \Phi; \mathbf{x}) &= \nabla_{\Phi} \mathbb{E}_{q_{\Phi}(z|\mathbf{x})} \left[ \log \frac{P_{\Theta}(\mathbf{x}, z)}{q_{\Phi}(z|\mathbf{x})} \right] \\ &= \nabla_{\Phi} \mathbb{E}_{q_{\Phi}(z|\mathbf{x})} [\log P_{\Theta}(\mathbf{x}, z) - \log q_{\Phi}(z|\mathbf{x})] \\ &\neq \mathbb{E}_{q_{\Phi}(z|\mathbf{x})} [\nabla_{\Phi} \log P_{\Theta}(\mathbf{x}, z) - \nabla_{\Phi} \log q_{\Phi}(z|\mathbf{x})]\end{aligned}$$

A reparametrization (variable change) trick allows a workaround

74

## Encoder function

Let us rewrite the decoder function with a random vector  $\epsilon$  whose distribution is not parametrized by  $\Phi$ :

$$\mathbf{z} = g(\mathbf{x}, \epsilon; \Phi)$$

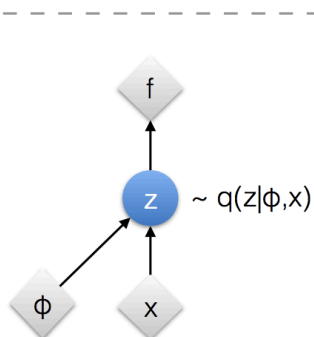
$$\begin{aligned} \nabla_{\Phi} \mathcal{L}(\Theta, \Phi; \mathbf{x}) &= \nabla_{\Phi} \mathbb{E}_{p(\epsilon)} \left[ \log \frac{P_{\Theta}(\mathbf{x}, \mathbf{z})}{q_{\Phi}(\mathbf{z}|\mathbf{x})} \right] \\ &= \nabla_{\Phi} \mathbb{E}_{p(\epsilon)} [\log P_{\Theta}(\mathbf{x}, \mathbf{z}) - \log q_{\Phi}(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_{p(\epsilon)} [\nabla_{\Phi} \log P_{\Theta}(\mathbf{x}, \mathbf{z}) - \nabla_{\Phi} \log q_{\Phi}(\mathbf{z}|\mathbf{x})] \\ &\approx \nabla_{\Phi} \log P_{\Theta}(\mathbf{x}, \mathbf{z}) - \nabla_{\Phi} \log q_{\Phi}(\mathbf{z}|\mathbf{x}) \end{aligned}$$

We just have to compute  $\log q_{\Phi}(\mathbf{z}|\mathbf{x})$  after the change of variable

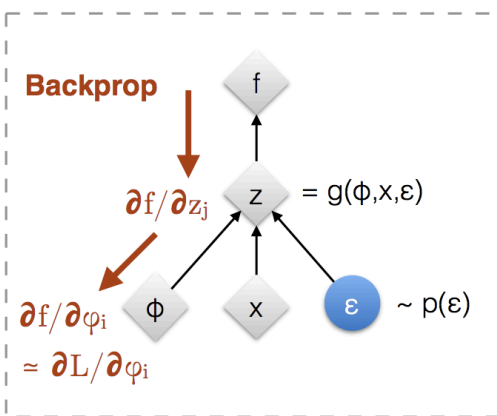
75

## Reparameterization trick

Original form



Reparameterised form



◆ : Deterministic node

● : Random node

[Kingma, 2013]  
 [Bengio, 2013]  
 [Kingma and Welling 2014]  
 [Rezende et al 2014]

76

## Computing $\log q_{\Phi}(\mathbf{z}|\mathbf{x})$ with a change of variable

$$\log q_{\Phi}(\mathbf{z}|\mathbf{x}) = \log p(\boldsymbol{\epsilon}) - \log d_{\Phi}(\mathbf{x}, \boldsymbol{\epsilon})$$

where the second term is the log of the absolute value of the determinant of the Jacobian matrix:

$$\log d_{\Phi}(\mathbf{x}, \boldsymbol{\epsilon}) = \log \left| \det \begin{pmatrix} \frac{\partial z_1}{\partial \epsilon_1} & \cdots & \frac{\partial z_1}{\partial \epsilon_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_k}{\partial \epsilon_1} & \cdots & \frac{\partial z_k}{\partial \epsilon_k} \end{pmatrix} \right|$$

variable change  $g()$  is chosen for the logdet being computationally affordable/simple

77

## Factorized Gaussian Posterior

### Model

$$(\boldsymbol{\mu}, \log \boldsymbol{\sigma}) = \text{EncoderNN}_{\Phi}(\mathbf{x})$$

$$q_{\Phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$$

$$q_{\Phi}(\mathbf{z}|\mathbf{x}) = \prod_i q_{\Phi}(z_i|\mathbf{x}) = \prod_i \mathcal{N}(z_i | \text{EncoderNN}_{\Phi}(\mathbf{x})) = \prod_i \mathcal{N}(z_i | \mu_i, \sigma_i^2)$$

### Reparametrization

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}$$

78

# Factorized Gaussian Posterior

The Jacobian of the transformation is

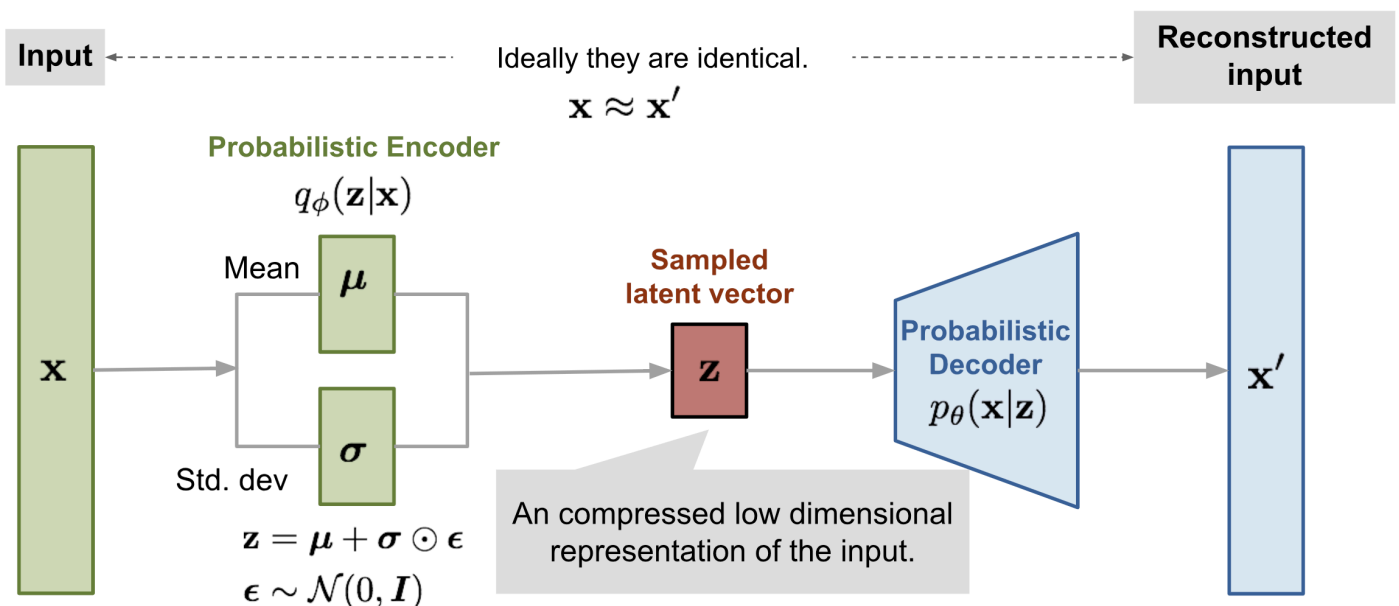
$$\log d_{\Phi}(\mathbf{x}, \boldsymbol{\epsilon}) = \log \left| \det \begin{pmatrix} \frac{\partial z_1}{\partial \epsilon_1} & \dots & \frac{\partial z_1}{\partial \epsilon_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_k}{\partial \epsilon_1} & \dots & \frac{\partial z_k}{\partial \epsilon_k} \end{pmatrix} \right| = \log \prod_i \sigma_i$$

The log posterior density is

$$\begin{aligned} \log q_{\Phi}(\mathbf{z}|\mathbf{x}) &= \log p(\boldsymbol{\epsilon}) - \log \left| \det \left( \frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} \right) \right| \\ &= \sum_i \log \mathcal{N}(\epsilon_i; 0, 1) - \log \sigma_i \end{aligned}$$

79

# Reparametrization trick



80

## Full Gaussian posterior

### Model

$$q_{\Phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

### Reparametrization

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\mathbf{z} = \boldsymbol{\mu} + \mathbf{L}\boldsymbol{\epsilon}$$

Where  $\mathbf{L}$  is a lower triangular matrix obtained from a Cholesky decomposition of  $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T$

81

## Full Gaussian posterior

The Jacobian has a simple form

$$\frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} = \mathbf{L}$$

As the determinant of a triangular matrix is the product of its diagonal terms,

$$\begin{aligned} \log q_{\Phi}(\mathbf{z}|\mathbf{x}) &= \log p(\boldsymbol{\epsilon}) - \log \left| \det\left(\frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}}\right) \right| \\ &= \sum_i \log \mathcal{N}(\epsilon_i; 0, 1) - \log L_{ii} \end{aligned}$$

82

## Algorithm input and output

Input:

- $\mathbf{X}$ : Dataset
- $h(\mathbf{z})$  decoding function, with dist.  $p_{\Theta}(\mathbf{x}, \mathbf{z})$
- $g(\mathbf{x})$  encoding function with dist.  $q_{\Psi}(\mathbf{z}|\mathbf{x})$

Output:

- $\Theta$
- $\Phi$

## Algorithm

Initialisation of  $\Theta$  and  $\Phi$

While SGD not converged do



- Draw a random minibatch  $\mathbf{X}^M \in \mathbf{X}$
- $\boldsymbol{\epsilon} \sim p(\boldsymbol{\epsilon})$  (Random noise for every datapoint in  $\mathbf{X}^M$ )
- Compute

$$\rightarrow \tilde{\mathcal{L}}_{\Theta, \Phi}(\mathbf{X}^M, \boldsymbol{\epsilon}) = \frac{1}{m} \sum_{\mathbf{x} \in \mathbf{X}^M} \left( \log p_{\Theta}(\mathbf{x}, \mathbf{z}) - \underbrace{\log q_{\Phi}(\mathbf{z}|\mathbf{x})}_{\log p(\boldsymbol{\epsilon}) - \log d_{\Phi}(\mathbf{x}, \boldsymbol{\epsilon})} \right) \text{ and}$$

→ its gradients

$$\rightarrow \nabla_{\Theta} \tilde{\mathcal{L}}_{\Theta, \Phi}(\mathbf{X}^M, \boldsymbol{\epsilon}) = \frac{1}{m} \sum_{\mathbf{x} \in \mathbf{X}^M} \frac{\partial \log p_{\Theta}(\mathbf{x}, \mathbf{z})}{\partial \Theta}$$

$$\rightarrow \nabla_{\Phi} \tilde{\mathcal{L}}_{\Theta, \Phi}(\mathbf{X}^M, \boldsymbol{\epsilon}) = \frac{1}{m} \sum_{\mathbf{x} \in \mathbf{X}^M} \frac{\log \partial d_{\Phi}(\mathbf{x}, \boldsymbol{\epsilon})}{\partial \Phi}$$

$$\bullet \begin{pmatrix} \Theta \\ \Phi \end{pmatrix} := \begin{pmatrix} \Theta \\ \Phi \end{pmatrix} + \eta \begin{pmatrix} \nabla_{\Theta} \tilde{\mathcal{L}}_{\Theta, \Phi}(\mathbf{X}^M, \boldsymbol{\epsilon}) \\ \nabla_{\Phi} \tilde{\mathcal{L}}_{\Theta, \Phi}(\mathbf{X}^M, \boldsymbol{\epsilon}) \end{pmatrix}$$

## Original example from Kingma: Gaussian model with MLP parametrization

### Multivariate Gaussian decoder with a diagonal covariance structure

Decoder  $p_{\Theta}(\mathbf{x}|\mathbf{z})$  or decoder (just swap  $\mathbf{x}$  and  $\mathbf{z}$ ) are assumed to have multivariate Gaussian dist. with a diagonal covariance structure:

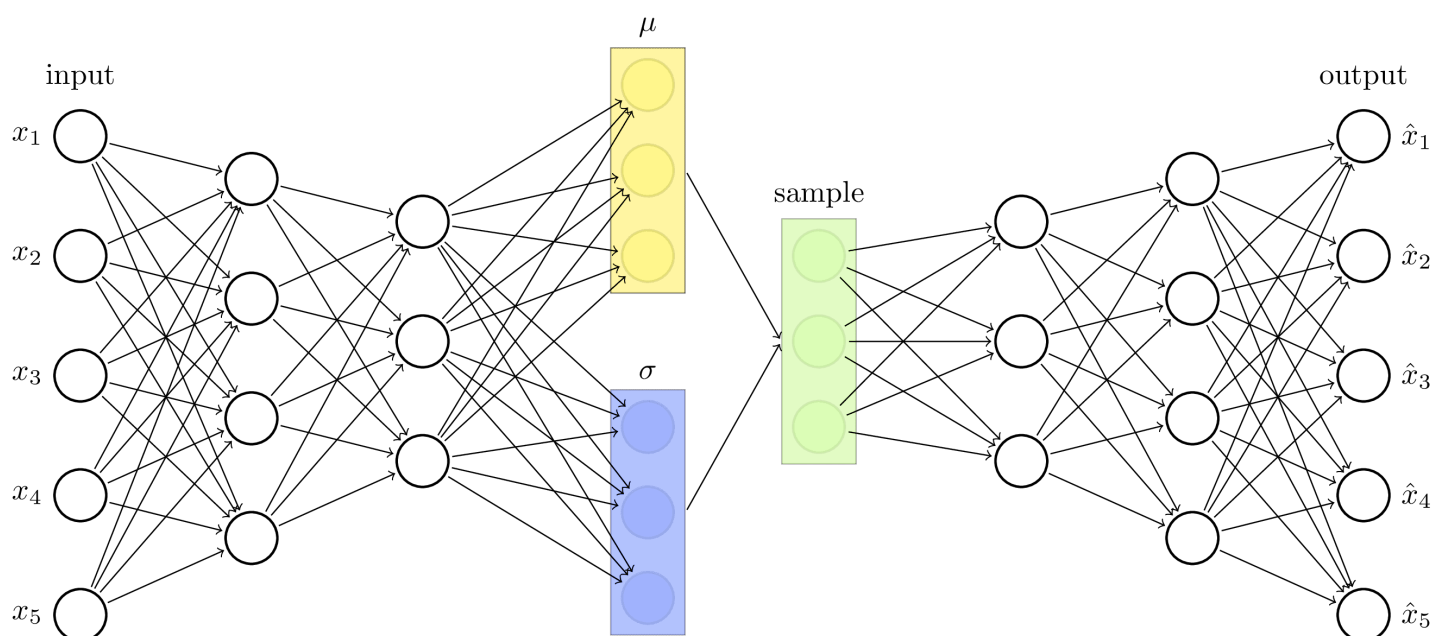
#### Decoder

- $\log p(\mathbf{x}|\mathbf{z}) = \log \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \sigma^2 I)$  where  $\boldsymbol{\mu} = W_4 \mathbf{h} + b_4$
- $\log \sigma^2 = W_5 \mathbf{h} + b_5$
- $\mathbf{h} = \tanh(W_3 \mathbf{z} + b_3)$

where  $\{W_3, W_4, W_5, b_3, b_4, b_5\}$  are the weights and biases of the MLP and part of  $\Theta$  when used as decoder.

- Let us consider the prior  $p_{\Theta}(\mathbf{z}) = \mathcal{N}(\mathbf{0}, I)$
- swap  $\mathbf{x}$  and  $\mathbf{z}$  in the decoder above to get  $q_{\Phi}(\mathbf{z}|\mathbf{x})$

## Gaussian VAE illustrated



# Mixture of Experts

87

## Historical Context

- Introduced in the early 1990s by Jacobs et al. (1991).

## *Ideas and Intuitions*

- **Inspired by cognitive science**, mimicking expert decision-making by assigning specialized sub-models to different tasks.
- **Enhances generalization** by combining multiple specialized models rather than relying on a single monolithic structure.
- **Divide-and-conquer approach** enables decomposition of complex tasks into manageable subtasks, improving prediction accuracy (Masoudnia & Ebrahimpour, 2014; Yuksel et al., 2012).
- **Widely applied in machine learning**, including speech recognition, image processing, LLM (Deepseek, Mixtral, ChatGPT...)

## Model Formulation

$$f(y_i; \mathbf{x}_i, \Theta) = \sum_{k=1}^K g_k(\mathbf{x}_i; \alpha_k) f_k(y_i; \mathbf{x}_i, \beta_k)$$

- $g_k(\mathbf{x}_i)$ : Output of the gating network

### Probabilistic formulation

- $f(y_i; \mathbf{x}_i, \Theta) = p(y_i | \mathbf{x}_i, \Theta)$
- the gating network comes from a latent variable  $z$ :

$$g_k(\mathbf{x}_i) = p(z_i = k | \mathbf{x}_i)$$

where

$$z_i = k | \mathbf{x}_i \sim M(1; g(\mathbf{x}_i) = (g_1(\mathbf{x}_i), \dots, g_K(\mathbf{x}_i)))$$

- $g$  is a sigmoid function for  $K = 2$  or softmax for  $K > 2$ .
- The expert  $f_k(y_i; \mathbf{x}_i, \beta_k)$  models a pdf corresponding the nature of  $y_i$

# MoE Architectures

- Different types of expert models (Gormley & Eisner, 2019).
- Variants include mixture models (unsupervised), prediction models (supervised).

## *Advantages & Challenges*

- MoEs enable large-scale models while reducing computational costs.
- Challenges: high parameter count, complex gating network training.

91

## Example of regression mixture

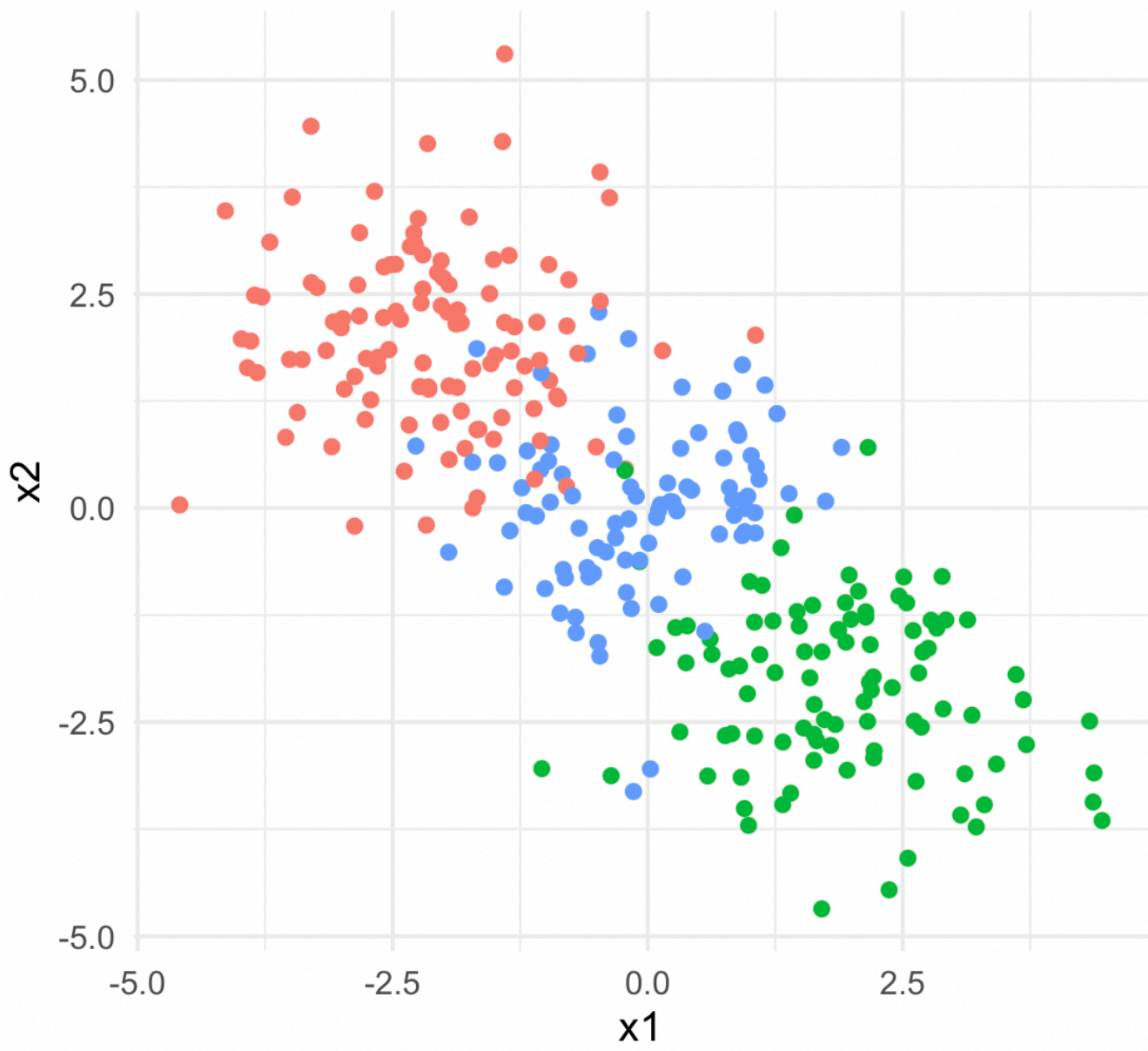
### *Training Process*

- **Labeled Data:**
  - The MoE is trained on a dataset where each input is associated with a known output.
- **Expert Specialization:**
  - Each expert learns to master a specific subset of the data, allowing for finer modeling of variations present in the data.
- **Gating Network Training:**
  - The gating network learns to associate inputs with the appropriate experts based on data characteristics.

### *Inference Phase*

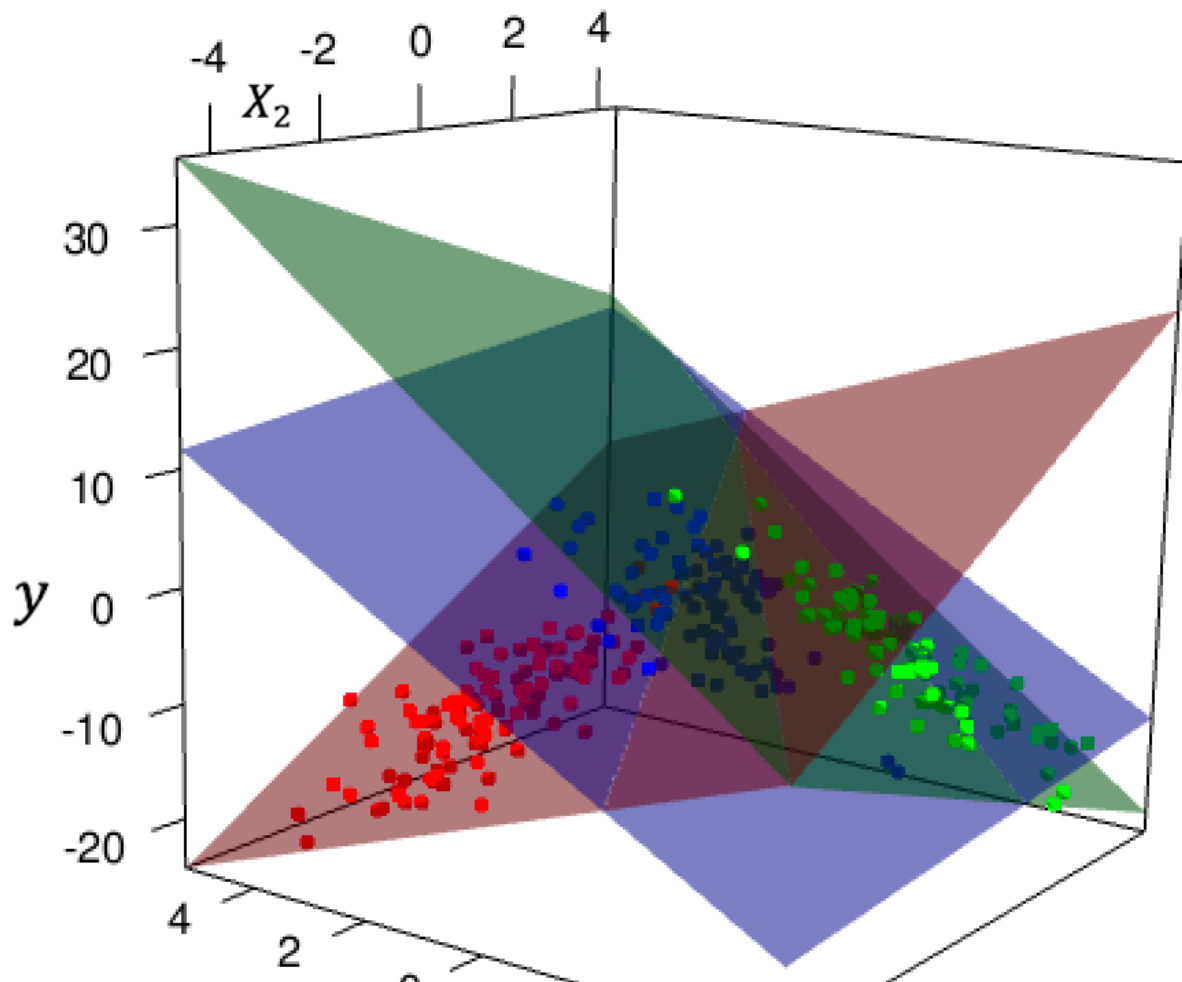
- **Input Routing:**
  - For new data, the gating network evaluates its characteristics and selects the most suitable expert to make the prediction.
- **Output Combination:**
  - In some cases, outputs from multiple experts may be combined to obtain a more robust final prediction.

## Example of regression mixture



Example of regression mixture

##



93

# Singular Value Decomposition

94



# Singular Value Decomposition

## Eigendecomposition of symmetric matrices

$\forall \mathbf{A} \in \mathbb{R}^{n \times n}$ , there exist an orthonormal matrix  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  and a diagonal matrix  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$

$$\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$$

## Singular Value Decomposition

Extend the decomposition to **rectangular matrices**

$$\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^T$$

95

# Applications in machine learning

- **Dimensionality Reduction:** SVD can be used for dimensionality reduction by reducing the rank of a matrix
- **Latent Semantic Analysis:** By decomposing a term-document matrix using SVD, LSA can capture the latent semantic structure of the data
- **Principal Component Analysis (PCA):** PCA is a SVD
- **Recommender Systems:** By factorizing the matrix using SVD, we can identify latent factors or features that capture underlying patterns and preferences.
- **Image Compression:** SVD is used in image compression techniques such as JPEG.
- **Matrix Completion:** SVD-based techniques are used in matrix completion problems, where missing or incomplete data needs to be imputed.
- ...

96

## Existence of the SVD for general matrices

For any matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , there exist two orthogonal matrices  $\mathbf{U} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{V} \in \mathbb{R}^{d \times d}$  and a nonnegative, "diagonal" matrix  $\mathbf{S} \in \mathbb{R}^{n \times d}$  such that

$$\mathbf{X}_{n \times d} = \mathbf{U}_{n \times n} \mathbf{S}_{n \times d} \mathbf{V}_{d \times d}^T$$

where  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$  and  $\mathbf{V}^T \mathbf{V} = \mathbf{I}$ .

### In a vector form

$$\mathbf{X}_{n \times d} = \sum_{j=1}^r S_{jj} \mathbf{u}_j \mathbf{v}_j^T$$

where  $r = \text{rank}(\mathbf{X})$ .

97

## Geometrical interpretation

Given any matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  it defines a linear transformation:

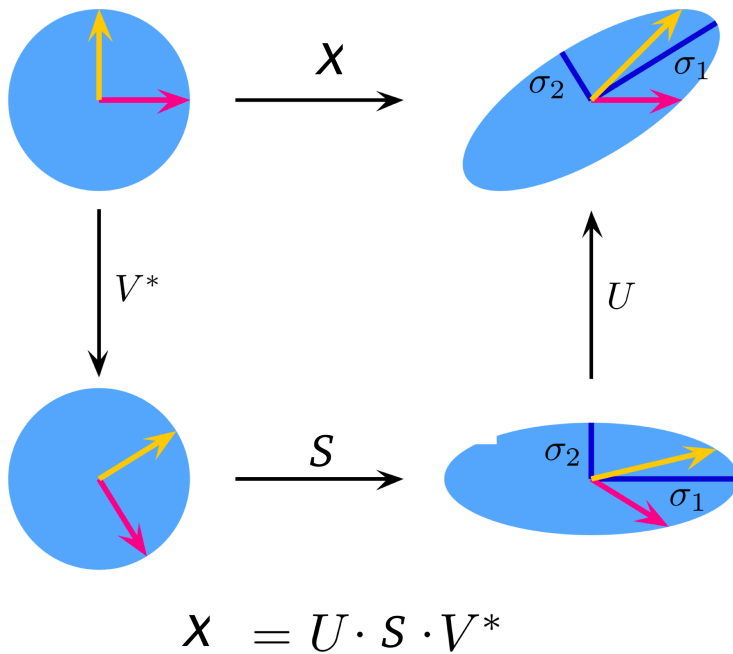
$$f: \mathbb{R}^d \rightarrow \mathbb{R}^n, f(\mathbf{x}) = \mathbf{X}\mathbf{x}.$$

The linear transformation  $f$  can be decomposed into three operations:

$$\underbrace{\mathbf{X}}_{\text{linear transformation}} \mathbf{x} = \underbrace{\mathbf{U}}_{\text{rotation}} \underbrace{\mathbf{S}}_{\text{scaling}} \underbrace{\mathbf{V}^T}_{\text{rotation}} \mathbf{x}$$

98

## Geometrical interpretation



99

## Different versions of SVD

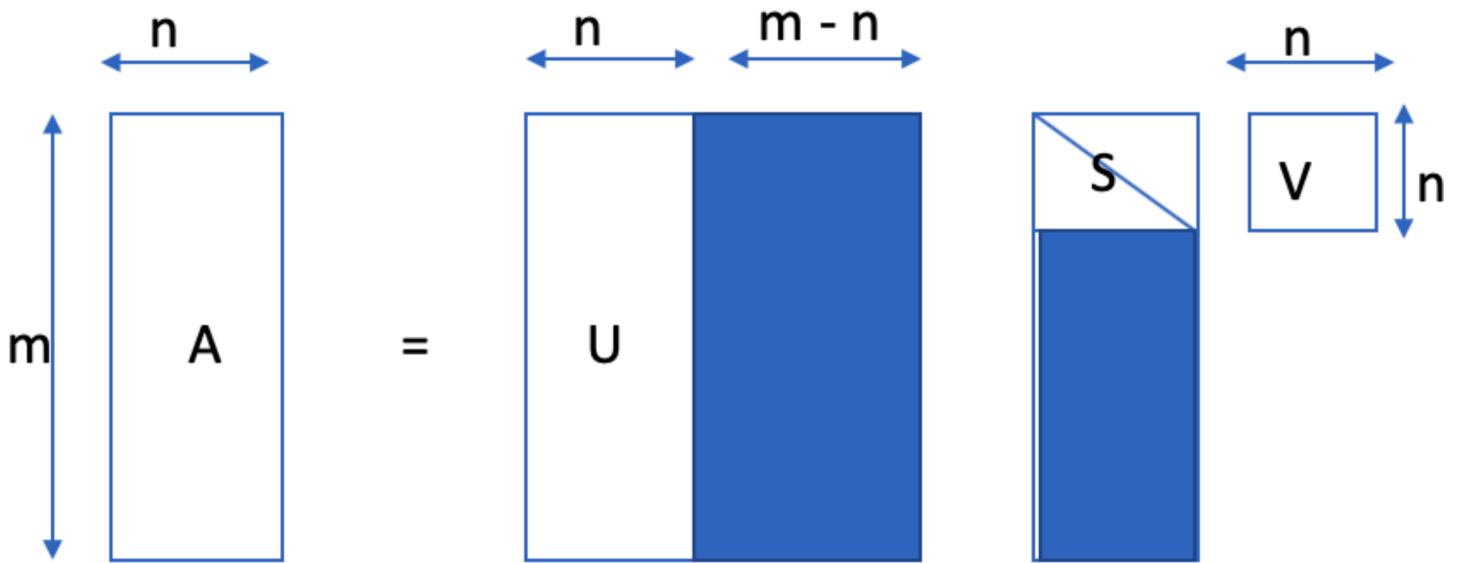
- Full SVD:

$$X_{n \times d} = U_{n \times n} S_{n \times d} V_{d \times d}^T$$

- Economy sized (thin, compact) SVD:

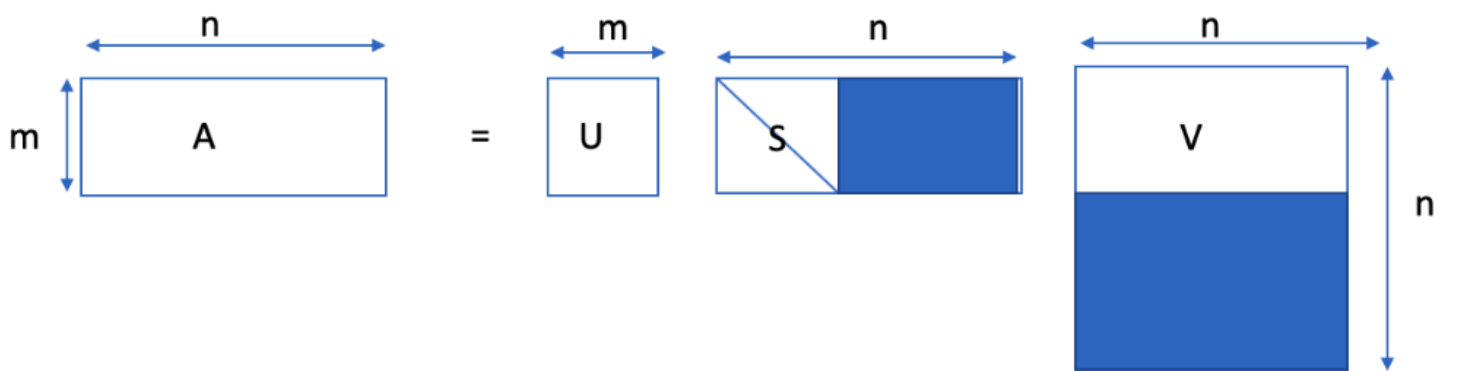
$$X_{n \times d} = U_{n \times r} S_{r \times r} V_{r \times d}^T$$

# SVD $n > d$



101

# SVD $n < d$



102

## Existence of the SVD

Consider  $\mathbf{A} = \mathbf{X}^T \mathbf{X} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T$  where  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_d)$  with  $\lambda_1 \geq \dots \geq \lambda_r > 0 = \lambda_{r+1} = \dots = \lambda_d$  (where  $r = \text{rank}(\mathbf{X}) \leq d$ ).

Let  $\sigma_i = \sqrt{\lambda_i}$  and correspondingly form the matrix

$$\mathbf{S}_{n \times d} = \begin{pmatrix} \text{diag}(\sigma_1, \dots, \sigma_r) & \mathbf{0}_{r \times (d-r)} \\ \mathbf{0}_{(n-r) \times r} & \mathbf{0}_{(n-r) \times (d-r)} \end{pmatrix}$$

Define also

$$\mathbf{u}_i = \frac{1}{\sigma_i} \mathbf{X} \mathbf{v}_i \in \mathbb{R}^n,$$

for each  $1 \leq i \leq r$ .

103

## Existence of the SVD

### Exercise

It is easy to show that the  $\mathbf{u}_1, \dots, \mathbf{u}_r$  are orthonormal vectors.

### Completion if needed

Choose  $\mathbf{u}_{r+1}, \dots, \mathbf{u}_n \in \mathbb{R}^n$  (through basis completion) such that

$$\mathbf{U} = [\mathbf{u}_1 \cdots \mathbf{u}_n] \in \mathbb{R}^{n \times n}$$

is an orthogonal matrix.

It verifies

$$\mathbf{X} \mathbf{V} = \mathbf{U} \mathbf{S},$$

i.e.,

104

## Existence of the SVD

$$\mathbf{X}[\mathbf{v}_1, \dots, \mathbf{v}_r, \mathbf{v}_{r+1}, \dots, \mathbf{v}_d] = [\mathbf{u}_1 \dots \mathbf{u}_r, \mathbf{u}_{r+1}, \dots, \mathbf{u}_n] \begin{pmatrix} \text{diag}(\sigma_1, \dots, \sigma_r) & \mathbf{0}_{r \times (d-r)} \\ \mathbf{0}_{(n-r) \times r} & \mathbf{0}_{(n-r) \times (d-r)} \end{pmatrix}$$

Two possible cases:

- $1 \leq i \leq r$  :  $\mathbf{X}\mathbf{v}_i = \sigma_i\mathbf{u}_i$  by construction.
- $i > r$  :  $\mathbf{X}\mathbf{v}_i = \mathbf{0}$ , which is due to  $\mathbf{X}^T\mathbf{X}\mathbf{v}_i = \mathbf{C}\mathbf{v}_i = \mathbf{0}\mathbf{v}_i = \mathbf{0}$ .

Consequently, we have obtained that

$$\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^T$$

105

## Properties

The linear application characterized by  $\mathbf{X}$  has the following properties:

- $\text{rank}(\mathbf{X}) = r$  is the number of non zero singular values
- $\text{kernel}(\mathbf{X}) = \text{span}(\mathbf{v}_{r+1}, \dots, \mathbf{v}_n)$
- $\text{range}(\mathbf{X}) = \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_r)$

106

# Low rank approximation of a matrix $X$

## Goal

Approximate a given matrix  $X$  with a rank- $k$  matrix, for a target rank  $k$ .

## Motivations

- Compression
- De-noising
- Matrix completion

107

## A first toy example

```
1 X<-matrix(c(1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12),4,3,byrow=TRUE)
2 X.svd<-svd(X)
3 cat("Original matrix:\n")
```

Original matrix:

```
1 print(X)
```

```
  [,1] [,2] [,3]
[1,]   1   2   3
[2,]   4   5   6
[3,]   7   8   9
[4,]  10  11  12
```

```
1 k<-2
2 cat("Approximation of rank 2:\n")
```

Approximation of rank 2:

```
1 print(X.svd$u[,1:k]**%diag(X.svd$d[1:k])**%t(X.svd$v[,1:k]))
```

```
  [,1] [,2] [,3]
[1,]   1   2   3
[2,]   4   5   6
[3,]   7   8   9
[4,]  10  11  12
```

```
1 cat("A basis of the column space:\n")
```

A basis of the column space:

```
1 print(X.svd$u[,1:k])
```

```
  [,1] [,2]
[1,] -0.1408767 -0.82471435
[2,] -0.3439463 -0.42626394
```

```
[3,] -0.5470159 -0.02781353  
[4,] -0.7500855  0.37063688
```

```
1 cat("\nA basis of the kernel:\n")
```

A basis of the kernel:

```
1 print(X.svd$u[,1:k])
```

```
      [,1]      [,2]  
[1,] -0.1408767 -0.82471435  
[2,] -0.3439463 -0.42626394  
[3,] -0.5470159 -0.02781353  
[4,] -0.7500855  0.37063688
```

108

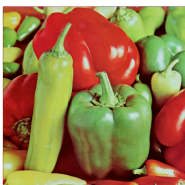
## Illustration of svd in image compression

See the demo of [Tim Baumann](#)

Example borrowed from [rich-d-wilkinson.github.io](http://rich-d-wilkinson.github.io)

The  $512 \times 512$  colour image is stored as three matrices R, B, G of the same dimension  $512 \times 512$  giving the intensity of red, green, and blue for each pixel. Naively storing this matrix requires 5.7Mb.

```
1 library(tiff)  
2 library(rasterImage)  
3 peppers<-readTIFF("../Silo-Images/Peppers.tiff")  
4 plot(as.raster(peppers))
```



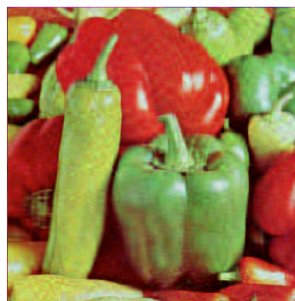
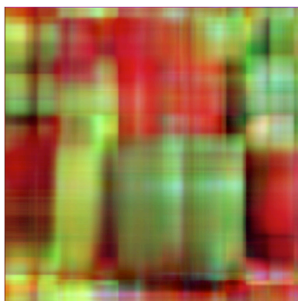
109



# Illustration of svd in image compression

Below the SVD of the three colour intensity matrices, and the view the image that results from using reduced rank versions with rank  $k \in \{5, 30, 100, 300\}$

```
1 svd_image <- function(im,k){
2   s <- svd(im)
3   Sigma_k <- diag(s$d[1:k])
4   U_k <- s$u[,1:k]
5   V_k <- s$v[,1:k]
6   im_k <- U_k %*% Sigma_k %*% t(V_k)
7   ## the reduced rank SVD produces some intensities <0 and >1.
8   # Let's truncate these
9   im_k[im_k>1]=1
10  im_k[im_k<0]=0
11  return(im_k)
12 }
13
14 par(mfrow=c(2,2), mar=c(1,1,1,1))
15
16 peprssvd<- peppers
17 for(k in c(4,30,100,300)){
18   svds<-list()
19   for(ii in 1:3) {
20     peprssvd[, ,ii]<-svd_image(peppers[, ,ii],k)
21   }
22   plot(as.raster(peprssvd))
23 }
```



## Low rank approximation of a matrix $\mathbf{X}$

### Frobenius norm

The Frobenius norm of a matrix  $\mathbf{X}$  is defined as

$$\|\mathbf{X}\|_F^2 = \sum_{ij} X_{ij}^2 = \text{trace}(\mathbf{X}^T \mathbf{X}) = \sum_{j=1}^r \sigma_j^2$$

### Rank $k$ matrix $\hat{\mathbf{X}}_k$

Let

$$\hat{\mathbf{X}}_k = \sum_{j=1}^k \sigma_j \mathbf{u}_j \mathbf{v}_j^T$$

111

## Low rank approximation of a matrix

For any matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  with non null singular values  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$

$$\hat{\mathbf{X}}_k = \arg \min_{\hat{\mathbf{X}}: \text{rank}(\hat{\mathbf{X}})=k} \|\mathbf{X} - \hat{\mathbf{X}}\|_F^2$$

$$\min_{\hat{\mathbf{X}}: \text{rank}(\hat{\mathbf{X}})=k} \|\mathbf{X} - \hat{\mathbf{X}}\|_F^2 = \sum_{j=k+1}^r \sigma_j^2$$

112

## Proof

We have

$$\|X - X_k\|_F^2 = \left\| \sum_{i=k+1}^n \sigma_i u_i v_i^\top \right\|_F^2 = \sum_{i=k+1}^n \sigma_i^2$$

We need to show that if  $Y_k = AB^\top$  where  $A$  and  $B$  have  $k$  columns then

$$\|X - X_k\|_F^2 = \sum_{i=k+1}^n \sigma_i^2 \leq \|X - Y_k\|_F^2.$$

113

## Proof

By the triangle inequality with the spectral norm, if  $X = X' + X''$  then  $\sigma_1(X) \leq \sigma_1(X') + \sigma_1(X'')$ .

Suppose  $X'_k$  and  $X''_k$  respectively denote the rank  $k$  approximation to  $X'$  and  $X''$  by SVD.

Then, for any  $i, j \geq 1$

$$\begin{aligned} \sigma_i(X') + \sigma_j(X'') &= \sigma_1(X' - X'_{i-1}) + \sigma_1(X'' - X''_{j-1}) \\ &\geq \sigma_1(X - X'_{i-1} - X''_{j-1}) \\ &\geq \sigma_1(X - X_{i+j-2}) \quad (\text{since } \text{rank}(X'_{i-1} + X''_{j-1}) \leq \text{rank}(X_{i+j-2})) \\ &= \sigma_{i+j-1}(X). \end{aligned}$$

114

## Proof

Since  $\sigma_{k+1}(Y_k) = 0$ , when  $X' = X - Y_k$  and  $X'' = Y_k$  we conclude that for  $i \geq 1, j = k + 1$

$$\sigma_i(X - Y_k) + \underbrace{\sigma_{k+1}(Y_k)}_0 \geq \sigma_{k+i}(X). \text{ Therefore,}$$

$$\|X - Y_k\|_F^2 = \sum_{i=1}^n \sigma_i(X - Y_k)^2 \geq \sum_{i=k+1}^n \sigma_i(X)^2 = \|X - X_k\|_F^2.$$

115

## Low rank approximation of a matrix and projection

If  $\text{rank}(\hat{X}) = k$ , then we can assume columns  $\hat{X}_i$  of  $\hat{X} \in E_k = \text{span}\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k\}$  where  $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k\}$  is a set of orthonormal vectors for the linear space of columns of  $X_k$ . First, observe that

$$\|X - \hat{X}\|_F^2 = \sum_i \|X_i - \hat{X}_i\|^2$$

### Optimum solution is the orthogonal projection

For each term  $\|X_i - v\|_2^2$ , the optimum solution is the projection of  $X_i$  onto  $E_k = \text{span}\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k\}$ :

$$\hat{X}_i = \sum_{j=1}^k \langle X_i, \mathbf{w}_j \rangle \mathbf{w}_j = \Pi_{E_k} X_i.$$

where  $\Pi_{E_k} = \sum_{j=1}^k \mathbf{w}_j \mathbf{w}_j^T$

116

## Projection on the orthogonal subspace

Consider  $\Pi_{E_k^\perp}$  the projection matrix on the space orthogonal to  $E_k$ . More precisely, let us add  $\mathbf{w}_{k+1}, \dots, \mathbf{w}_n$  such that  $\mathbf{w}_1, \dots, \mathbf{w}_n$  form an orthonormal basis of  $\mathbb{R}^n$ . Then,

$$\Pi_{E_k^\perp} = \sum_{j=k+1}^n \mathbf{w}_j \mathbf{w}_j^T$$

$$\|X - \hat{X}\|_F^2 = \|X - \Pi_{E_k} X\|_F^2 = \|(I - \Pi_{E_k})X\|_F^2 = \|\Pi_{E_k^\perp} X\|_F^2$$

117

## Relation to principal component analysis

### Warning

**$X$  is considered as centered.** This transformation (cloud translation allows considerable simplification)

### Decomposition of $X$

Considering the orthogonal projection on  $E_k$

$$X = \Pi_{E_k} X + \Pi_{E_k^\perp} X$$

118

## Criterion

$$\|\mathbf{X}\|_F^2 = \underbrace{\|\Pi_{E_k} \mathbf{X}\|_F^2}_{\text{approximation}} + \underbrace{\|\Pi_{E_k^\perp} \mathbf{X}\|_F^2}_{\text{error}}$$

In terms of inertia, PCA maximizes the projected inertia (approximation) while minimizing the distances to the space of projection (error):

$$I_T = I_E + I_{E_k^\perp}$$

119

## Best low rank approximation

$$\hat{\mathbf{X}}_k = \Pi_{E_k} \mathbf{X} = \sum_{j=1}^k \sigma_j \mathbf{u}_j \mathbf{v}_j^T = \mathbf{U}_{\bullet,1:k} \mathbf{S}_{1:k,1:k} \mathbf{V}_{\bullet,1:k}^T$$

where  $\mathbf{X} = \underbrace{\mathbf{U}}_{n \times n} \underbrace{\mathbf{S}}_{n \times d} \underbrace{\mathbf{V}^T}_{d \times d}$

120

## Different views of the approximation

The approximation

$$\|\mathbf{X} - \hat{\mathbf{X}}\|_F^2$$

can be considered in multiple ways:

- approximation of the row
- approximation of the columns

### Notations

If  $\mathbf{X}$  is a data table,

- each row  $\mathbf{x}_i^T$  is a description of an individual
- each column  $\mathbf{X}_j$  is variable describing  $n$  individuals

121

## Rows approximation (projection of the individuals)

Transposing the matrix the best low rank approximation becomes

$$\hat{\mathbf{X}}^T_k = \Pi_{F_k} \mathbf{X}^T = \sum_{j=1}^k \sigma_j \mathbf{v}_j \mathbf{u}_j^T = \mathbf{V}_{\bullet, 1:k} \mathbf{S}_{1:k, 1:k} \mathbf{U}_{\bullet, 1:k}^T$$

where  $F_k = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$

122

## The approximation error

$$\|\mathbf{X} - \hat{\mathbf{X}}\|_F^2 = \|\mathbf{X}^T - \hat{\mathbf{X}}^T\|_F^2 = \sum_i \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|_2^2$$

Each row  $\mathbf{x}_i$  is approximated by

$$\hat{\mathbf{x}}_i = \Pi_{F_k} \mathbf{x}_i = V_{F_k} V_{F_k}^T \mathbf{x}_i$$

where  $V_{F_k}$  is the matrix composed of the vectors defining  $F_k$ .

123

## Projection of the variables

$\Pi_{E_k}$  is the projection matrix on  $E = \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$

$$\Pi_E = U_{E_k} U_{E_k}^T$$

and

$$\Pi_{E_k} \mathbf{X} = U_{1:n,1:k} \underbrace{U_{1:n,1:k}^T U_{1:n,1:n}}_{(I_k, \mathbf{0}_{k,n-k})} \mathbf{S}_{1:n,1:k} V_{1:d,1:d}^T = U_{1:n,1:k} \mathbf{S}_{1:k,1:k} V_{1:d,1:k}^T$$

124



## k first principal components

$$\mathbf{C}_{1:n,1:k} = \mathbf{U}_{E_k} \mathbf{S}_{1:k,1:k}$$

where  $\mathbf{S}_{1:k,1:k} = \text{diag}(\sigma_1, \dots, \sigma_k)$ .

The principal component are the coordinates of the projection of the rows of  $\mathbf{X}$  on  $F_k$ :

$$\mathbf{C}_{1:n,1:k} = \mathbf{XV}_{1:d,1:k}$$

125

## Percentage of information

We have  $\mathbf{C}_{\bullet,1:k}^T \mathbf{C}_{\bullet,1:k} = \mathbf{S}_{1:k,1:k}^2 = \text{diag}(\sigma_1^2, \dots, \sigma_k^2)$ , thus

$$\|\mathbf{C}_{\bullet,1:k}\|_F^2 = \sum_{j=1}^k \sigma_j^2$$

and

$$\frac{\|\mathbf{C}_{\bullet,1:k}\|_F^2}{\|\mathbf{X}\|_F^2} = \frac{\sum_{j=1}^k \sigma_j^2}{\sum_{j=1}^d \sigma_j^2} \in [0, 1]$$

126

## Correlations

$$\widehat{\text{cor}}(\mathbf{X}_{\bullet,j}, \mathbf{C}_{\bullet,k}) = \frac{\mathbf{X}_{\bullet,j}^T \mathbf{C}_{\bullet,k}}{\|\mathbf{X}_{\bullet,j}\| \|\mathbf{C}_{\bullet,k}\|} = \cos(\widehat{\mathbf{X}_{\bullet,j}}, \widehat{\mathbf{C}_{\bullet,k}})$$

127

## Duality

It is easy to show that

- the columns of  $\mathbf{V}$  are the eigenvector of  $\mathbf{X}^T \mathbf{X}$
- the columns of  $\mathbf{U}$  are the eigenvector of  $\mathbf{X} \mathbf{X}^T$

Thus the principal component of  $\mathbf{X}^T \mathbf{X}$  are the eigenvectors of  $\mathbf{X} \mathbf{X}^T$  and vice-versa

128

# Multi Dimensional Scaling

129

## Dimensionality reduction and manifold learning

### Goal

Given pairwise

- dissimilarities  $\delta_{ij}$  or
- high-dimensional input data  $\{\mathbf{x}_i\}_{i=1 \dots n}$

reconstruct a lower dimensional map with embedded data  $\{\mathbf{y}_i\}_{i=1 \dots n}$

### Principle

- Minimize an objective function quantifying the discrepancies between the  $\delta_{ij}$  and the distances  $d(\mathbf{y}_i, \mathbf{y}_j)$
- w.r.t  $\{\mathbf{y}_i\}_{i=1 \dots n}$ .

## A Brief history

The methodology of Multidimensional Positioning (Multidimensional Scaling, MDS) was born in the USA in the 1950s

- Works of Torgerson (1952, 1958),
- Works of Shepard (1962),

initial applications to sparse data in

- psychometrics,
- marketing,
- sensory analysis.

131

## Proximity measures

### *Distance*

The function  $d$  from  $E \times E$  in  $\mathbb{R}^+$  is a distance if it satisfies the following properties:

- $\forall i, j \in E, d_{ij} = d_{ji}$
- $\forall i \in E, d_{ii} = 0$
- $\forall i, j \text{ and } k \in E, d_{ik} \leq d_{ij} + d_{jk}$

### *Dissimilarity*

The function  $\delta$  of  $E \times E$  in  $\mathbb{R}^+$  is a distance if it satisfies the following properties:

- $\forall i, j \in E, \delta_{ij} = \delta_{ji}$
- $\forall i \in E, \delta_{ii} = 0$

132

# Classical scaling

Also known as Torgerson scaling, principal coordinate analysis (PCoA), ...

## Principle

- assume that the dissimilarities are distances and
- find the main coordinates that explain the distances.

133

## From distances to scalar product

- How to calculate the distance matrix from  $X$ ?

$$\begin{aligned}d_{ij}^2 &= \sum_{a=1}^p (x_{ia}^2 + x_{ja}^2 - 2x_{ia}x_{ja}) \\ &= \sum_{a=1}^p x_{ia}^2 + \sum_{a=1}^p x_{ja}^2 - 2 \sum_{a=1}^p x_{ia}x_{ja}\end{aligned}$$

from where

$$D^2 = \text{diag}(XX^t)\mathbb{I}_{(1,n)} + \mathbb{I}_{(n,1)}\text{diag}(XX^t)^t - 2XX^t$$

## Distance and data table (continued)}

- Let  $J$  be the centering matrix,

$$J = I - \frac{1}{n} \mathbb{I}_{(n,n)}$$

- Centering of  $X$  in columns

$$JX = X - \frac{1}{n} \mathbb{I}_{(n,n)} X$$

- Centering of  $X$  in rows

$$X^t J = X^t - \frac{1}{n} X^t \mathbb{I}_{(n,n)} X$$

## Distance and data table (continued)}

*Problem: we know  $D^2$ , we want  $X$ : Double centering*

$$\begin{aligned} -\frac{1}{2} J D^2 J &= -\frac{1}{2} J \text{diag}(X X^t) * \mathbb{I}_{(1,n)} J \\ &\quad - \frac{1}{2} J \mathbb{I}_{(n,1)} \text{diag}(X X^t)^t J \\ &\quad + J X X^t J \\ &= X X^t \end{aligned}$$

# Low rank approximation and principal coordinates

The idea is to minimize

$$\|XX^T - YY^T\|_F^2$$

w.r.t  $\text{rank}(YY^T) = k$

We know the minimum is a low rank approximation  $B = YY^T$  of rank  $k$

$$B = (U_{E_k} S_{1:k,1:k}) (S_{1:k,1:k}^T U_{E_k}^T) = YY^T$$

- the analysis of the triple consists in replacing the matrix  $D^2$  by the matrix of the square of dissimilarities  $\Delta^2$

## Remark

If  $\Delta^2$  is a distance matrix the solution  $Y = U_{E_k} S_{1:k,1:d}$  (first  $k$  principal components of  $X^T X$ ) is optimal and PCoA is equivalent to PCA.

136

## Algorithm

1. Compute the matrix  $\Delta^2$
2. Double centering of  $\Delta^2$ :

$$B_{\Delta^2} = -\frac{1}{2} J \Delta^2 J$$

3. Spectral decomposition of  $B_{\Delta^2}$ :

$$B_{\Delta^2} = U \Lambda U^t$$

4. Let  $k$  be the representation dimension chosen for the solution.

- $\Lambda_+$  is the matrix of  $k$  largest eigenvalues, listed in descending order on the diagonal.
- $U^+$  is the matrix of  $k$  corresponding eigenvectors
- **The solution of the problem** is

$$Y = U^+ \Lambda_+^{\frac{1}{2}}$$

137

## Optimized criterion

- the optimized criterion is

$$\begin{aligned}L(X) &= \left\| -\frac{1}{2}J(D^2(X) - \Delta^2)J \right\|^2 \\ &= \left\| XX^t + \frac{1}{2}J\Delta^2J \right\|^2 \\ &= \left\| XX^t - B_{\Delta^2} \right\|^2\end{aligned}$$

with  $B_{\Delta^2}$  of rank  $k$

- Remarks: If  $\Delta$  is a dissimilarity matrix then some principal components will be negative!

138

## Models and functions for the MDS

MDS matches proximities  $\delta_{ij}$  to distances  $d_{ij}$  between objects:

$$f : \delta_{ij} \longrightarrow d_{ij}(Y)$$

with  $Y$  a classical array of low dimensional data - the dissimilarities  $\delta_{ij}$  are the data of the pb - the distances  $d_{ij}(Y)$  are the unknowns

- the MDS finds a configuration  $Y$  in a space of dimension  $m$ , used to calculate distances  $d_{ij}$

139



## In practice

the dissimilarities are marred by errors and one does not look for  $f$  such that

$$f(\delta_{ij}) = d_{ij}(Y)$$

but rather

$$f(\delta_{ij}) \approx d_{ij}(Y)$$

140

## Error function

### Definition of an error function

$$e_{ij} = (f(\delta_{ij}) - d_{ij}(Y))^2$$

### Raw global error (raw stress)

$$\sigma_r^2(Y) = \sum_{i>j} (f(\delta_{ij}) - d_{ij}(Y))^2$$

- The global error is not very informative~: A large value of  $\sigma_r^2(Y)$  does not necessarily indicate bad result.

141

## Normalized error function

### Normalized Stress

$$\sigma_n^2(Y) = \frac{\sum_{i>j} \left( f(\delta_{ij}) - d_{ij}(Y) \right)^2}{\sum_{i>j} d_{ij}^2(Y)}$$

### A general form

$$\sigma^2(Y) = \sum_{i>j} w_{ij} \left( f(\delta_{ij}) - d_{ij}(Y) \right)^2$$

142

## Different approaches

- metric: quantitative approach
- non-metric : qualitative approach

### Metric approach

Algebraic transformation of  $\delta_{ij}$  (e.g.  $f(\delta_{ij}) = a\delta_{ij} + b$ )

With

$$f(\delta_{ij}) = \delta_{ij} + b(1 - \mathbb{I}_{(i=j)})$$

assurance of the existence of an  $Y$  configuration of dimension  $n - 2$  at most such that

$$f(\delta_{ij}) = d_{ij}(Y)$$

### Non-metric approach

143

# Sammon's Projection

Metric Method: "Projection" into a space of low dimension (1,2, or 3)

- $Y = (\mathbf{y}_1, \dots, \mathbf{y}_n)^t$  configuration sought
- Euclidian distance

$$\begin{aligned}d_{ij} &= d(\mathbf{y}_i, \mathbf{y}_j) \\ &= \|\mathbf{y}_i - \mathbf{y}_j\|_2\end{aligned}$$

144

## Sammon stress function

Sammon searches for the  $\mathbf{y}_i$  minimizing

$$S(Y) = \frac{1}{\sum_{i < j} \delta_{ij}} \frac{\sum_{i < j} (\delta_{ij} - d_{ij}(Y))^2}{\delta_{ij}}$$

- $f$  identity
- $w_{kl} = \frac{1}{(\sum_{i < j} \delta_{ij}) \delta_{kl}}$

### Remarks:

Sammon's Stress

- is invariant to rotations, translation and scaling
- focuses on small distances

145

# Minimisation of Sammon's Stress

- minimisation of a function from  $\mathbb{R}^{n \times k}$  to  $\mathbb{R}$
- Global minimum is difficult to reach
- Gradient descent

## Remarks

- Many possible optimisation methods
- Initial Random configuration

146

# Example in R

## *sammon* from library MASS

Sammon's Non-Linear Mapping

### Description

One form of non-metric multidimensional scaling.

### Usage

```
sammon(d, y = cmdscale(d, k), k = 2, niter = 100, trace = TRUE,  
magic = 0.2, tol = 1e-4)
```

147

## Sammon gradient

Let show that

$$\frac{\partial S(Y)}{\partial y_{ik}} = -2 \sum_{j < i} w_{ij} \frac{\delta_{ij} - d_{ij}}{d_{ij}} (y_{ik} - y_{jk})$$

From chain rule

$$\frac{\partial S(Y)}{\partial y_{ik}} = \underbrace{\frac{\partial S(Y)}{\partial \|y_i - y_j\|}}_1 \times \underbrace{\frac{\partial \|y_i - y_j\|}{\partial y_{ik}}}_2$$

1.  $\frac{\partial S(Y)}{\partial d_{ij}} = -2 \sum_{j < i} w_{ij} (\delta_{ij} - d_{ij})$
2.  $\frac{\partial \|y_i - y_j\|}{\partial y_{ik}} = \frac{\partial d_{ij}}{\partial d_{ij}^2} \frac{\partial d_{ij}^2}{\partial y_{ik}} = \frac{y_{ik} - y_{jk}}{d_{ij}}$

148

## Sammon gradient descent

At iteration  $q$

$$y_{ik}^{(q+1)} := y_{ik}^{(q)} - \eta \frac{\partial S(Y)}{\partial y_{ik}}$$

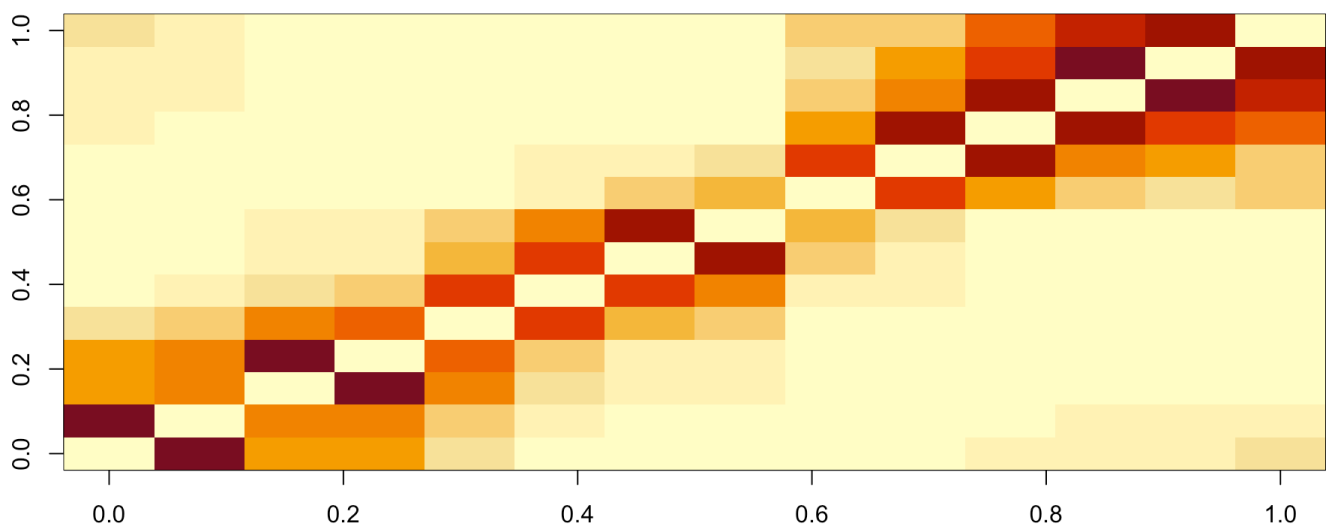
149

## Example Ekman colors

Ekman presents similarities for 14 colors which are based on a rating by 31 subjects where each pair of colors was rated on a 5-point scale (0 = no similarity up to 4 = identical). After averaging, the similarities were divided by 4 such that they are within the unit interval. Similarities of colors with wavelengths from 434 to 674 nm.

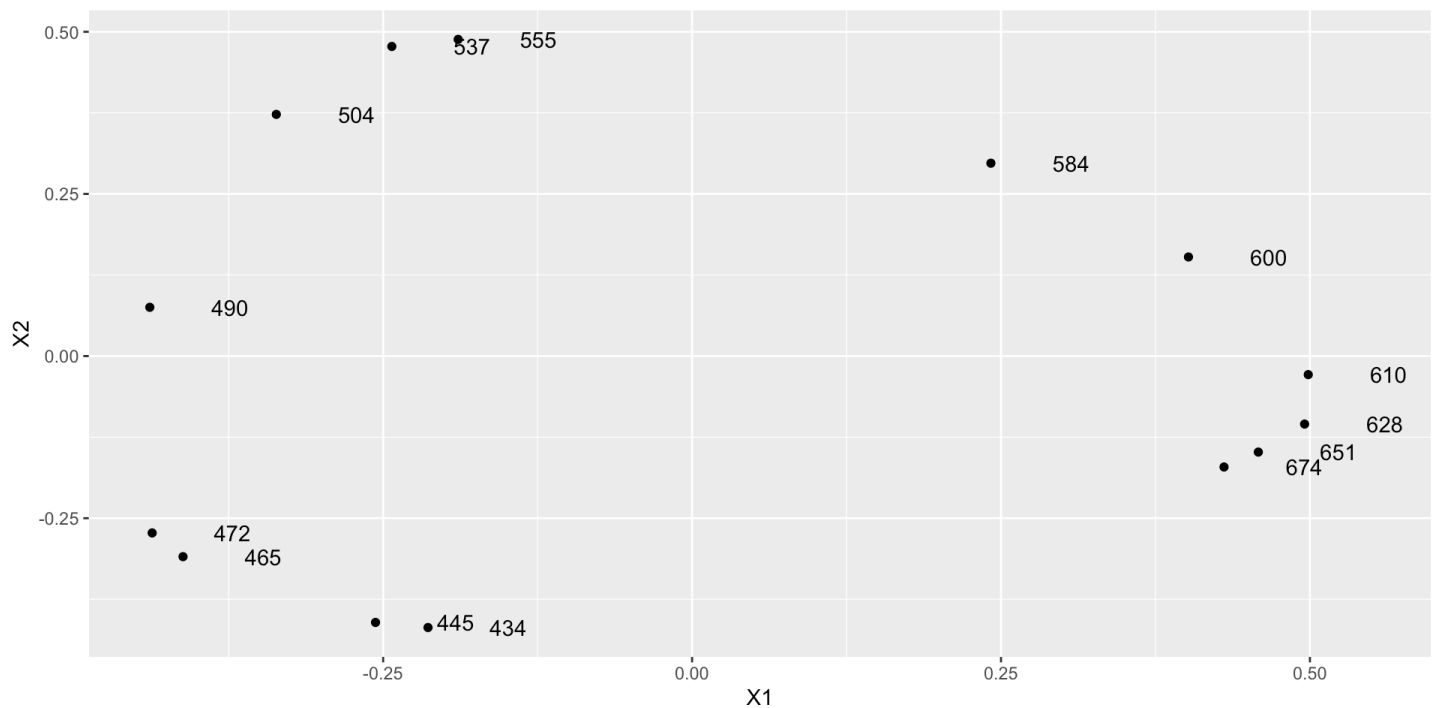
150

## Example Ekman colors



151

## Example Classical scaling



152

## t-SNE

Variation of Stochastic Neighbor Embedding

- easier to optimize,
- produces significantly better visualizations by reducing the tendency
- scalable: for very large data sets, t-SNE can use random walks on neighborhood graphs

153

## Stochastic Neighbor Embedding

Stochastic Neighbor Embedding (SNE) converts distances into conditional probabilities that represent similarities

*Conditional probabilities  $p_{j|i}$*

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / (2\sigma_i^2))}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / (2\sigma_i^2))}.$$

*Interpretation*

Probability that  $x_i$  would pick  $x_j$  as its neighbor if neighbors were picked in proportion to their probability density under a Gaussian centered at  $x_i$ :

154

## Stochastic Neighbor Embedding

similar conditional probability between the  $y_i$ , which we denote by  $q_{j|i}$ .

*Conditional probabilities,  $q_{j|i}$*

$$q_{j|i} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq i} \exp(-\|y_i - y_k\|^2)}.$$

with  $q_{i|i} = 0$ .

*SNE can also be applied directly to similarity data*

provided similarities can be interpreted as conditional probabilities we set

155



## Cost function

A natural measure of the faithfulness with which  $q_{j|i}$  models  $p_{j|i}$  is the Kullback-Leibler divergence (equal to the cross-entropy up to an additive constant).

$$C = \sum_i KL(P_i || Q_i) = \sum_i \sum_j p_{j|i} \log \frac{p_{j|i}}{q_{j|i}},$$

the SNE cost function focuses on retaining the local structure of the data in the map (for *reasonable values of the variance* of the Gaussian in the high-dimensional space).

156

## Selecting the variance $\sigma_i$

SNE performs a binary search for the value of  $\sigma_i$  that produces a  $P_i$  with a fixed perplexity that is specified by the user

### Perplexity

$$Perp(P_i) = 2^{H(P_i)},$$

where  $H(P_i)$  is the Shannon entropy of  $P_i$  measured in bits

$$H(P_i) = - \sum_j p_{j|i} \log_2 p_{j|i}.$$

### Remarks

- the perplexity increases monotonically with the variance  $\sigma_i$
- smooth measure of the effective number of neighbors
- Typical values are between 5 and 50.

157

# Gradient descent

## Gradient

$$\frac{\partial C}{\partial \mathbf{y}_i} = 2 \sum_j (\mathbf{p}_{j|i} - \mathbf{q}_{j|i} + \mathbf{p}_{i|j} - \mathbf{q}_{i|j})(\mathbf{y}_i - \mathbf{y}_j).$$

## Interpretation

The gradient may be interpreted as the resultant force created by a set of springs between the map point  $\mathbf{y}_i$  and all other map points  $\mathbf{y}_j$ .

- force exerted by the spring between  $\mathbf{y}_i$  and  $\mathbf{y}_j$  is proportional to its length
- also proportional to its stiffness  $(\mathbf{p}_{j|i} - \mathbf{q}_{j|i} + \mathbf{p}_{i|j} + \mathbf{q}_{i|j})$

## Initialisation

$n$  points randomly sampled from an isotropic Gaussian with small variance that is centered around the origin

158

# Gradient in practice

- To speed up the optimization and to avoid poor local minima, a relatively large momentum term is added to the gradient:

$$\mathbf{y}^{(t)} = \mathbf{y}^{(t-1)} + \eta \frac{\partial C}{\partial \mathbf{y}} + \alpha(t) (\mathbf{y}^{(t-1)} - \mathbf{y}^{(t-2)}),$$

where  $\mathbf{y}^{(t)}$  indicates the solution at iteration  $t$ ,  $\eta$  indicates the learning rate, and  $\alpha(t)$  represents the momentum at iteration  $t$ .

- In addition, in the early stages of the optimization, Gaussian noise is added to the map points after each iteration. Gradually reducing the variance of this noise performs a type of simulated annealing

159

# t-Distributed Stochastic Neighbor Embedding

- cost function of SNE is difficult to optimize
- t-SNE suffers from the “crowding problem”

## t-SNE vs SNE

1. Uses a symmetrized version of the SNE cost function with simpler gradients
2. Student-t distribution rather than a Gaussian to compute the similarity between two points *in the low dimensional space*.

160

## Symmetric SNE

$$C = KL(P||Q) = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}},$$

where  $p_{ii}$  and  $q_{ii}$  to zero.

- $p_{ij} = p_{ji}$  and
- $q_{ij} = q_{ji}$  for all  $i, j$ .

### Joint probabilities $p_{ij}$

$$p_{ij} = \frac{\exp(-\|x_i - x_j\|^2/2\sigma^2)}{\sum_{k \neq l} \exp(-\|x_k - x_l\|^2/2\sigma^2)}$$

### Joint probabilities $q_{ij}$

$$q_{ij} = \frac{\exp(-\|x_i - x_j\|^2)}{\sum_{k \neq l} \exp(-\|x_k - x_l\|^2)}.$$

161

## t-SNE

Problems when a high-dimensional datapoint  $x_i$  is an outlier, the values of  $p_{ij}$  are extremely small for all  $j$ , so the location of its low-dimensional map point  $y_i$  has very little effect on the cost function.

- $p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$ . ensures that  $\sum_j p_{ij} > \frac{1}{2n} \forall x_i$ ,
- each  $x_i$  makes a significant contribution to the cost function.

### t-SNE Gradient

The main advantage of the symmetric version of SNE is the simpler form of its gradient, which is faster to compute.

$$\frac{\partial C}{\partial y_i} = 4 \sum_j (p_{ij} - q_{ij})(y_i - y_j).$$

symmetric SNE produces maps that are just as good as asymmetric SNE

162

## The Crowding Problem

the area of the two-dimensional map that is available to accommodate moderately distant datapoints will not be nearly large enough compared with the area available to accommodate nearby datapoints.

- Modeling the small distances accurately has detrimental effect on moderate distances representation
- the crowding problem is not specific to SNE

163

## Mismatched tails can compensate for mismatched dimensionalities

### *t*-distribution for the $q_{ij}$

Student distribution with a single degree of freedom also known as a Cauchy distribution

$$q_{ij} = \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}}$$

### Gradient

$$\frac{\partial C}{\partial \mathbf{y}_i} = 4 \sum_j (p_{ij} - q_{ij})(\mathbf{y}_i - \mathbf{y}_j)(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}. \quad (1)$$

## Two optimization tricks

### “Early compression”

- Force the map points to stay close together at the start of the optimization.
- Easy for clusters to move through one another and thus to explore the space of solutions
- implemented by adding an additional L2-penalty to the cost function that is proportional

### “Early exaggeration”

Multiply all of the  $p_{ij}$ 's by, for example, 4, in the initial stages of the optimization.

- natural clusters in the data tend to form tight widely separated clusters in the map.
- creates a lot of relatively empty space in the map
- makes it much easier for the clusters to move around relative to one another

165

## Uniform manifold approximation and projection (UMAP)

UMAP works directly with similarities and uses cross-entropy as a criterion.

### Distance in high dimensional space

$$v_{j|i} = \exp [(-d(x_i, x_j) - \rho_i)/\sigma_i]$$

- $d(x_i, x_j)$  is the distance between  $x_i$  and  $x_j$ , which UMAP does not require to be Euclidean.
- $\rho_i$  is the distance to the nearest neighbor of  $i$ .
- $\sigma_i$  is the normalizing factor, which is chosen by a specific Algorithm and plays a similar role to the perplexity-based calibration of  $\sigma_i$  in t-SNE.

Symmetrization is carried out by fuzzy set union

$$v_{ij} = v_{j|i} + v_{i|j} - v_{j|i}v_{i|j}$$

166

# Uniform manifold approximation and projection (UMAP)

## Low dimensional similarities

$$w_{ij} = (1 + a\|y_i - y_j\|_2^{2b})^{-1}$$

where  $a$  and  $b$  are user-defined positive values

## Cost function (cross-entropy)

$$C_{UMAP} = \sum_{i \neq j} v_{ij} \log \frac{v_{ij}}{w_{ij}} + (1 - v_{ij}) \log \frac{1 - v_{ij}}{1 - w_{ij}}$$

Optimization process is done by stochastic gradient descent

167

## UMAP vs t-SNE

In terms of performance,

- UMAP is often considered to produce more stable and consistent results than t-SNE,
- especially when dealing with large datasets.

On the other hand,

- t-SNE is often considered to produce more visually appealing results,
- especially when dealing with small datasets.

168

# Word Embedding

169

## Words and Vectors

Vector or distributional models of meaning are generally based on a co-occurrence matrix

### *Two common co-occurrence matrix*

- the word-document matrix
- the word-word matrix.



## Word-document matrix

In a Word-document matrix, each row represents a word in the vocabulary and each column represents a document from some collection of documents.

	<b>As You Like It</b>	<b>Twelfth Night</b>	<b>Julius Caesar</b>	<b>Henry V</b>
<b>battle</b>	1	0	7	13
<b>good</b>	114	80	62	89
<b>fool</b>	36	58	1	4
<b>wit</b>	20	15	2	3

*Number of words occurrence in 4 Shakespeare plays (from speech and language processing)*

171

## Word-Word matrix

In the Word-context matrix the columns are labeled by words rather than documents.

- each cell records the number of times the row (target) word and the column (context) word co-occur in some context in some training corpus.

### Context

- a context could be the document, in which case the cell represents the number of times the two words appear in the same document.
- smaller contexts are common: generally a window around the word

	<b>aardvark</b>	...	<b>computer</b>	<b>data</b>	<b>result</b>	<b>pie</b>	<b>sugar</b>	...
<b>cherry</b>	0	...	2	8	9	442	25	...
<b>strawberry</b>	0	...	0	0	1	60	19	...
<b>digital</b>	0	...	1670	1683	85	5	4	...
<b>information</b>	0	...	3325	3982	378	5	13	...

word-word matrix from a Wikipedia corpus (from speech and language processing)

172

## Cosine for measuring similarity

The cosine similarity metric (or empirical correlation) between two word vectors  $\mathbf{v}$  and  $\mathbf{w}$  (lines of a word-document or word-word matrix)

$$\cos(\mathbf{v}, \mathbf{w}) = \frac{\mathbf{v}^T \mathbf{w}}{\|\mathbf{v}\| \|\mathbf{w}\|}$$

173

## TF-IDF: Weighing terms in the vector

### Term frequency (tf)

- Raw frequency  $tf_{t,d} = \text{count}(t, d)$  is very skewed and not very discriminative.
- Usually squashed by using the  $\log_{10}$ : a word appearing 100 times in a document doesn't make that word 100 times more relevant

$$tf_{t,d} = \log_{10} (\text{count}(t, d) + 1)$$

### Inverse document frequency

The second factor in tf-idf is used to give a higher weight to words that occur only in a few documents.

$$idf_t = \log_{10} \left( \frac{N}{df_t} \right)$$

where  $N$  is the total number of documents, and  $df_t$  is the number of documents in which term  $t$  occurs.

174

## TF-IDF: Weighing terms in the vector

The tf-idf weighted value  $w_{t,d}$  for word  $t$  in document  $d$  thus combines term frequency  $tf_{t,d}$  with  $idf$ :

$$w_{t,d} = tf_{t,d} \times idf_t$$

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	0.074	0	0.22	0.28
good	0	0	0	0
fool	0.019	0.021	0.0036	0.0083
wit	0.049	0.044	0.018	0.022

A tf-idf weighted term-document matrix. Note that the tf-idf weighting has eliminated the importance of the ubiquitous word *good* and vastly reduced the impact of the almost-ubiquitous word *fool*.

175

## Pointwise Mutual Information (PMI)

An alternative weighting function to tf-idf, PPMI (positive pointwise mutual information), is used for **term(target)-term(context)-matrices**.

- Target word, word we focus on.
- Context words surrounding the target word

The pointwise mutual information between a target word  $w$  and a context word  $c$  (Church and Hanks 1989, Church and Hanks 1990) is defined as:

$$PMI(w, c) = \log_2 \frac{P(w, c)}{P(w)P(c)}$$

176

## Positive PMI (called PPMI)

Negative PMI values (which imply things are co-occurring less often than we would expect by chance) tend to be unreliable unless our corpora are enormous.

### *Positive PMI*

replaces all negative PMI values with zero (Church and Hanks 1989, Dagan et al. 1993, Niwa and Nitta 1994)

$$PPMI(w, c) = \max\left(\log_2 \frac{P(w, c)}{P(w)P(c)}, 0\right)$$

177

## Applications of the tf-idf or PPMI vector models

### Principle

- Computing two documents (or word) similarity after transformation.
- Given two documents  $d_1$  and  $d_2$ , the similarity is  $\cos(d_1, d_2)$ .

### Applications

- **Documents** : information retrieval, plagiarism detection, news recommender systems, and even for digital humanities tasks like comparing different versions of a text to see which are similar to each other.
- **Words** : finding word paraphrases, tracking changes in word meaning, or automatically discovering meanings of words in different corpora.

178

## Latent Semantic Analysis (LSA)

### Goal and assumptions

- extracting and representing the underlying meaning of words in a corpus of texts.
- words that occur in similar contexts have similar meanings.

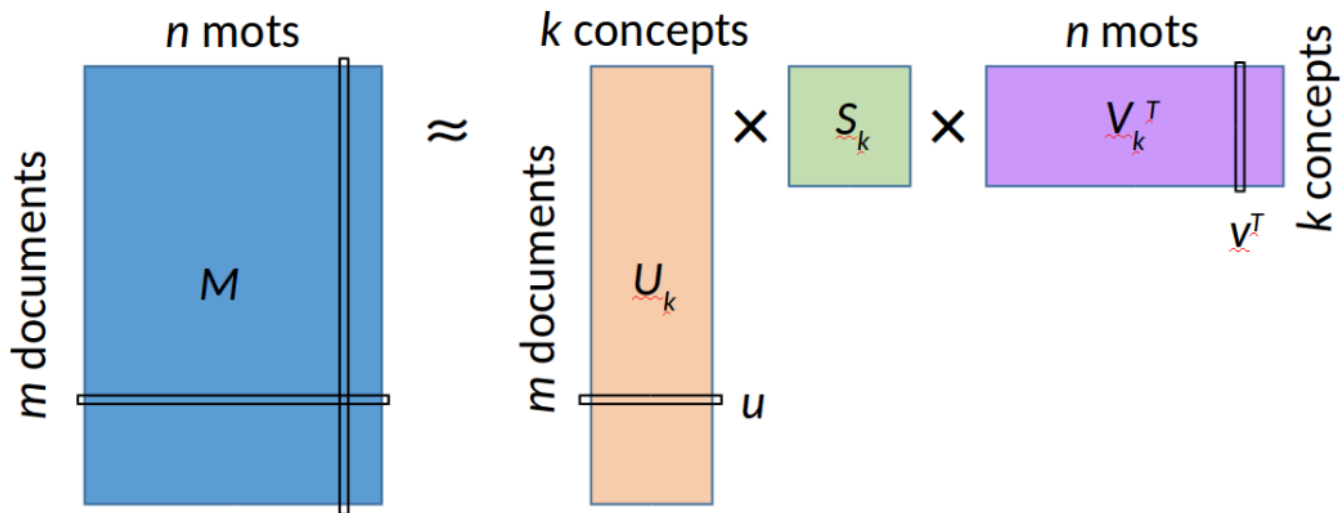
### Principle

LSA uses singular value decomposition (SVD) to identify latent, or hidden, patterns in word co-occurrence data.

179

# Latent Semantic Analysis

LSA uses SVD on the matrix of word-document matrix to identify the latent concepts (contexts, topics, ...) that underlie the relationships between words in the corpus:  $W = USV^T$  with classical low rank approximation  $U_k S_k V_k^T$



180

# Latent Semantic Analysis

- The values on the main diagonal of  $S_k$  indicate the 'importance' of each of the  $k$  main 'latent concepts' (or factors).
- For each of the document, the corresponding row of  $U_k$  allows us to see which concepts are present and with what weights.
- For each of the concept, the associated column of  $V_k$  indicates which terms form the concept (and with what weights).

Calculating the similarity between a term and a document involves choosing:

- the row  $\mathbf{u}$  corresponding to the document in the matrix  $U_k$ ,
- the row  $\mathbf{v}$  corresponding to the term in the matrix  $V_k$ ,
- and then computing the product  $\mathbf{u}S_k\mathbf{v}^T$ .

To determine the similarity of a term to each of the documents

$$U_k S_k v^T$$

Conversely, we can determine the similarities between a document  $\mathbf{u}$  and all the terms through

$$\mathbf{u} \mathbf{S}_k \mathbf{V}_k^T$$

181

## Latent Dirichlet Allocation (LDA)

### *History*

In the context of population genetics, LDA was proposed by J. K. Pritchard, M. Stephens and P. Donnelly in 2000.

LDA was applied in machine learning by David Blei, Andrew Ng and Michael I. Jordan in 2003.

### *Principle*

Latent Dirichlet Allocation (LDA) is a probabilistic generative model that assumes that

- each document is generated by a mixture of latent topics
- each topic is generated by a mixture of words from the vocabulary.

Documents are represented as random mixtures over latent topics, where each topic is characterized by a distribution over all the words.

182

## LDA Generative process

For a corpus  $D$  consisting of  $M$  documents

1. Choose topic parameter vector for document  $i$ :  $\theta_i \sim \text{Dir}(\alpha)$ , where  $i \in \{1, \dots, M\}$  and  $\text{Dir}(\alpha)$  is a Dirichlet distribution with a symmetric parameter  $\alpha$  which typically is sparse ( $\alpha$ ).
2. Choose the word parameters vector of topic  $k$ :  $\varphi_k \sim \text{Dir}(\beta)$ , where  $k \in \{1, \dots, K\}$  and  $\beta$  typically is sparse
3. For each of the word positions  $i, j$ , where  $i \in \{1, \dots, M\}$ , and  $j \in \{1, \dots, N_i\}$ 
  - a. Choose a topic  $z_{i,j} \sim \text{Cat}(\theta_i)$ .
  - b. Choose a word  $w_{i,j} \sim \text{Cat}(\varphi_{z_{i,j}})$ .

183

## Latent Dirichlet Allocation (LDA) in summary

$$\varphi_{k=1 \dots K} \sim \text{Dirichlet}_V(\beta)$$

$$\theta_{d=1 \dots M} \sim \text{Dirichlet}_K(\alpha)$$

$$z_{d=1 \dots M, w=1 \dots N_d} \sim \text{Categorical}_K(\theta_d)$$

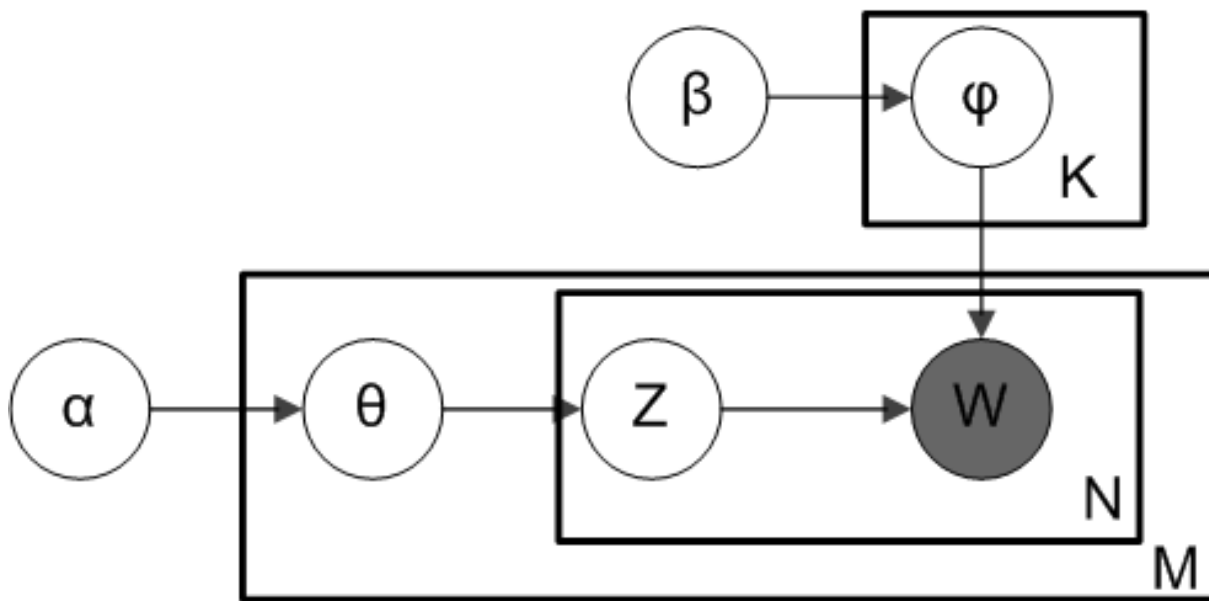
$$w_{d=1 \dots M, w=1 \dots N_d} \sim \text{Categorical}_V(\varphi_{z_{dw}})$$

$$P(\mathbf{W}, \mathbf{Z}, \boldsymbol{\theta}, \boldsymbol{\varphi}; \alpha, \beta) = \prod_{i=1}^K P(\varphi_i; \beta) \prod_{j=1}^M P(\theta_j; \alpha) \prod_{t=1}^N P(Z_{j,t} | \theta_j) P(W_{j,t} | \varphi_{Z_{j,t}}),$$

184



## Latent Dirichlet Allocation (LDA) estimation



Estimating the parameters could be achieved through Variational Bayes EM algorithm.

185

## Example of LDA in Action

Corpus of Three Documents

1. "cats dogs pets love"
2. "dogs bark loud outside"
3. "politics government election law"

Step 1: Identify Topics

After running LDA, it may discover two topics (dist. over words): - Topic 1 (Pets): {cats, dogs, pets, bark, love} - Topic 2 (Politics): {politics, government, election, law}

Step 2: Assign Topics to Documents

- Document 1 -> 90% Topic 1, 10% Topic 2
- Document 2 -> 80% Topic 1, 20% Topic 2
- Document 3 -> 95% Topic 2, 5% Topic 1

186

## Pro and Cons

### *Advantages of LDA*

- Unsupervised Learning: No labeled data needed.
- Interpretable Topics: Extracts meaningful topics from text.

### *Limitations of LDA*

- Fixed Number of Topics: Must specify the number of topics beforehand.
- Bag-of-Words Assumption: Ignores word order and context.
- Computationally Expensive: Requires approximation methods for inference.

### *Applications of LDA*

187

## Word2vec

### *Embeddings*

- more powerful word representation
- short dense vectors:  $d$  ranging from 50-1000, rather than the much larger vocabulary size  $|V|$  or number of documents  $D$
- dense vectors work better in every NLP task than sparse vectors

### *Skip-gram with negative sampling*

- The skip-gram algorithm is one of two algorithms in a software package called word2vec
- The other algorithm is CBOW (continuous bag of words)

188

## Embedding derived from classification

*Data : words and their context (neighboring words ( $\pm L$ ))*

... lemon, a [tablespoon of apricot jam, a] pinch ... c1 c2 w c3 c4

### *Classification task*

Given a tuple  $(\mathbf{w}, \mathbf{c}) \in \mathbb{R}^d \times \mathbb{R}^d$  of a target word  $\mathbf{w}$  paired with a candidate context word  $\mathbf{c}$  (for example (apricot, jam), or perhaps (apricot, aardvark)) return the probability that  $\mathbf{c}$  is a real context word (true for jam, false for aardvark):

$$P(y_{wc} = 1 | \mathbf{w}, \mathbf{c}) = 1 - P(y_{wc} = 0 | \mathbf{w}, \mathbf{c})$$

### *Similarity*

we rely on the intuition that two vectors are similar if they have a high dot product (after all, cosine is just a normalized dot product). In other words:

189

## From dot product to logistic regression

### *Logit*

$$P(y_{wc} = 1 | \mathbf{w}, \mathbf{c}) = \frac{1}{1 + \exp(-\mathbf{c}^T \mathbf{w})}$$

### *Criterion*

The criterion is a kind of likelihood including positive examples (observed association) and negative examples (non observed associations). In fact skip-gram with negative sampling (SGNS) uses more negative examples than positive examples (with the ratio between them set by a parameter  $k$ ).

$$L = \sum_{\mathbf{w}, \mathbf{c}: y_{wc}=1} \log P(y_{wc} = 1 | \mathbf{w}, \mathbf{c}) + \sum_{\mathbf{w}', \mathbf{c}': y_{w'c'}=0} \log P(y_{w'c'} = 0 | \mathbf{w}', \mathbf{c}')$$

190

## Example

- The learning algorithm for skip-gram embeddings takes as input a corpus of text
- random embedding vector initialization for each of the N vocabulary words
- iteratively shift the embedding of each word  $w$  to be more like the context words

Let's start by considering a single piece of training data:

... lemon, a [tablespoon of apricot jam, a] pinch ...

### positive examples +

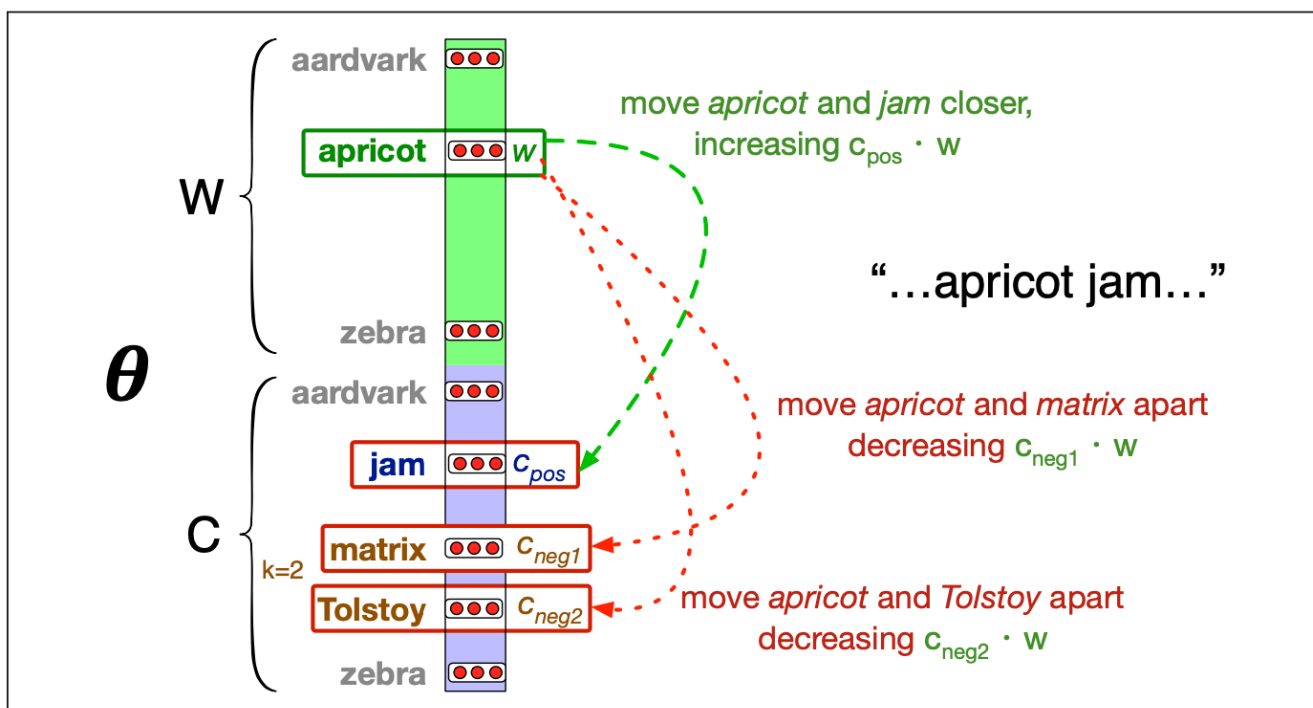
$w$	$c_{pos}$
apricot	tablespoon
apricot	of
apricot	jam
apricot	a

### negative examples -

$w$	$c_{neg}$	$w$	$c_{neg}$
apricot	aardvark	apricot	seven
apricot	my	apricot	forever
apricot	where	apricot	dear
apricot	coaxial	apricot	if

191

## Skip-gram model learns two separate embeddings



192

## Final embeddings

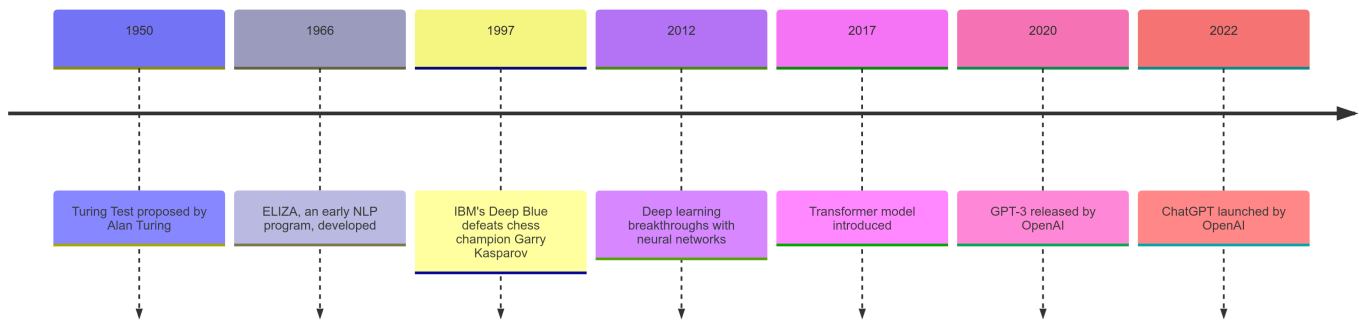
- skip-gram outputs the target matrix  $W$  and the context matrix  $C$ .
- It's common to just add them together, representing word  $i$  with the vector  $w_i + c_i$ .
- Alternatively we can throw away the  $C$  matrix and just represent each word  $i$  by the vector  $w_i$ .
- As with the simple count-based methods like tf-idf, the context window size  $L$  affects the performance of skip-gram embeddings, and experiments often tune the parameter  $L$  on a devset.

# Transformers

# The nature of transformers

- Deep learning architecture based on **attention mechanisms**
- Weight the importance of different tokens in a sequence to **model long-range dependencies efficiently**

## History of AI Leading to ChatGPT



History of AI leading to chatGPT

# The nature of transformers

A classical transformer is a function

$$T(w_{1:n}) = \begin{bmatrix} p_1 \\ \vdots \\ p_{|V|} \end{bmatrix}$$

$$\mathcal{V} = \{A, \text{CAT}, \text{DOG}, \text{BARKS}, \text{MEOWS}\}$$

$$T \left[ w_{1:2} = \begin{matrix} A & \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \end{matrix} \right] = \begin{bmatrix} P(A | w_{1:2}) = 0 \\ P(\text{CAT} | w_{1:2}) = 0 \\ P(\text{DOG} | w_{1:2}) = 0 \\ P(\text{BARKS} | w_{1:2}) = 0.9 \\ P(\text{MEOWS} | w_{1:2}) = 0.1 \end{bmatrix}$$

The nature of transformer

- $\mathcal{V}$  is a finite set (vocabulary) (Open AI uses 50,000 token as vocabulary)
- $w_{1:n}$  be a sequence of token, where each  $w_i \in \mathcal{V}$  (In practice context varies from 1 to n, the blocksize)
- $p_j$  is the probability of  $w_{(j)}$  the jth token in  $\mathcal{V}$

## Remarks

## Optimized criteria

Transformers typically optimize different criteria depending on the task they are trained on.

### Causal Language Modeling (CLM) / Autoregressive Loss

- **Used in:** GPT models
- **Objective:** Predict the next token given the previous tokens

The training set is a set of sequences  $(s_j)$  of  $n$  tokens  $s_j = (w_1^j, w_2^j, \dots, w_n^j)$ :

$$L_{CLM} = - \sum_j \sum_{t=1}^n \sum_{c \in V} \mathbb{I}_{w_t^j=c} \log P(w_t^j = c | w_{1:t-1}^j)$$

where  $w_{1:t-1}^j$  denotes the tokens preceding position  $t$  in the sequence  $j$

The criterion is a kind of pseudo log-likelihood.

If nothing is learned each token has the same probability of appearing: the largest possible value of the criterion is

$$- \log \frac{1}{|V|}$$

## Masked Language Modeling (MLM)

- **Used in:** BERT models
- **Objective:** Predict randomly masked tokens in a sentence.

Let  $M_j \subset \{1, \dots, n\}$  be the set of masked positions in the sequence  $(s_j)$ :

$$L_{MLM} = - \sum_j \sum_{t \in M_j} \sum_{c \in V} \mathbb{I}_{w_t^j = c} \log P(w_t^j = c | w_{(1:n) \setminus M}^j)$$

198

## Conditional Language Modeling / Translation Loss

- **Used in:** T5, BART
- **Objective:** Predict an output sequence given an input sequence (e.g., translation, summarization).

Given input sequence  $w^j = (w_1^j, \dots, w_n^j)$  and output sequence  $v^j = (v_1^j, \dots, v_m^j)$ :

$$L_{Seq2Seq} = - \sum_j \sum_{t=1}^m \sum_{c \in V} \mathbb{I}_{v_t^j = c} \log P(v_m^j = c | v_{1:m-1}^j, w_{1:n}^j)$$

### Other losses

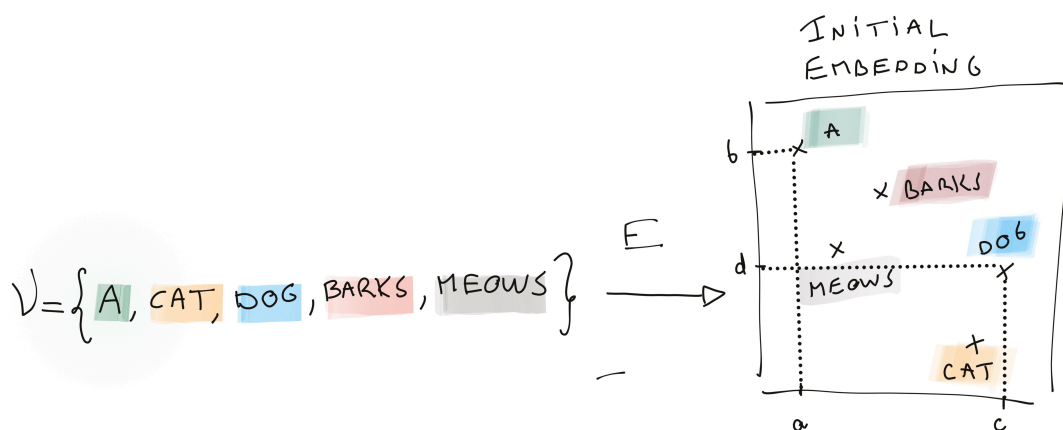
Transformers exact objective functions depends on type of considered problems

199



## Initial embedding

- The tokens are first embedded in a  $d$  dimensional space (see previous topic)
- the representation  $x_i = Ew_i$  is independent of context.
- Let  $X = x_{1:n} = w_{1:n}E^T \in \mathbb{R}^{n \times d}$ .



Initial embedding

200

## A minimal self-attention architecture for decoder

**Attention is a communication mechanism** between the tokens of a sequence.

Self attention proposes a contextual representation  $h_i$  of  $x_i$  as a linear combination the sequence:

$$h_i = V \cdot \sum_{j=1}^n \alpha_{ij} x_j$$

- $h_i$  aggregates the past between query (current position  $i$ ) and key (other past position  $j$ )  
→  $\sum_j \alpha_{ij} = 1$
- In decoder the future cannot communicate with the past  
→  $\forall j > i, \alpha_{ij} = 0$
- the structure of the attention matrix is thus lower triangular in classical GPT decoder,
- it is not necessary to have a triangular attention matrix (i.e. sentiment analysis)
- $V$  is learned matrix of parameters and can be seen a projection matrix modifying the embedding of the linear combination of sequence tokens. Let consider  $V$  as a square matrix  $d \times d$ .

201

## Attention and matrix notation

- $\alpha_{ij} = \text{softmax}(x_i^T Q^T K x_j) = \frac{\exp(x_i^T Q^T K x_j)}{\sum_{j=1}^n \exp(x_i^T Q^T K x_j)}$ , where

→ Q is a matrix that modifies the embedding of the token we are looking for. Let consider V as a square matrix  $d \times d$ .

↪ K is a matrix that modifies the embedding of the word we are comparing against. Let consider V as a square matrix  $d \times d$ .

### Matrix notation

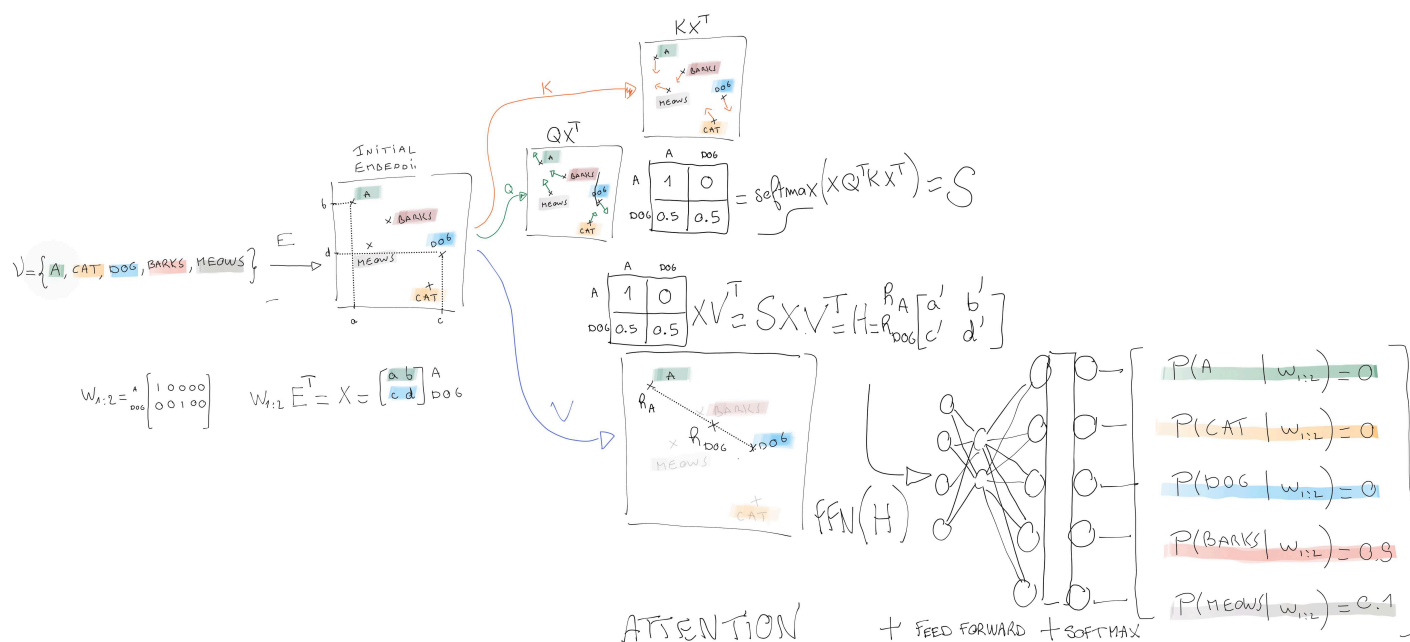
$$H = \begin{bmatrix} h_1^T \\ \vdots \\ h_n^T \end{bmatrix} = \text{Attention}(Q, K, V, X) = \text{softmax}(X(Q^T K)X^T)XV^T$$

### Remark

In the original paper ("Attention is all you need", 2017) the notation is different and notes:

- $XQ^T$  as  $Q$ ,
- $XK^T$  as  $K$ , and
- $XV^T$  as  $V$

# Illustration of self-Attention



Self-attention block in auto-encoder

## Interpretation of Q, K, V

- The query is modified version of an initial embedding  $q_i = Qx_i$
- The Key is modified version of an initial embedding  $k_i = Kx_i$
- Value is modified version of an initial embedding  $v_i = Vx_i$ , which should include contextual information since it is a transformation of a linear combination of the context tokens...

### Attention Concept

### Sometimes Similar to

Query

Target (What to focus on)

Key

Context (What to compare against, the source of information)

Value

Contextual embedding

## Important details

The original paper introduces many add-ons, which much improve the transformers performances:

- **Positional encoding** for introducing the notion of order in the sequence
- FFN: **Feed Forward Neural layer**.
  - As the attention is a linear operation composed with a softmax, a Feed Forward Neural layer is added for making the function more flexible
- **Multi-blocks**: In the spirit of Deep learning blocks of attention and FNN are repeated...
- **Dropout layers** are used for prevenring overfitting (from “Dropout a simple way to prevent NN from overfitting” Hinton)
- Simple trick of **Skip connection** for stabilizing the gradient (from “Deep learning residuals for image recognition”)
  - The deep learning structure tends to cause problems with gradient computation (vanishing or exploding gradient)
- In order to make tokens comparable, they are all **normalized** at the output of each blocks:
  - The classical z-transformed is applied (with Bayesian correction)
- **Multi-head attention**: computes multiple attention functions (heads) in parallel

205

## Positional encoding

- Attention do not consider order of tokens !
- Incorporating information about the order of tokens is achieved by adding positional encodings to the input embeddings

The original transformer implementation uses:

- even indices:

$$PE_{\text{pos},2i} = \sin\left(\frac{\text{pos}}{10000^{\frac{2i}{d}}}\right)$$

- odd indices:

$$PE_{\text{pos},2i+1} = \cos\left(\frac{\text{pos}}{10000^{\frac{2i}{d}}}\right)$$

It ensures that each position is assigned a unique encoding

206

## Feed Forward Neural Network

- **Attention mechanism operates linearly** to capture dependencies between tokens in a sequence
- Non-linearity enhance the model's expressive power

### Structure of the FFN:

Each FFN consists of two linear transformations with a non-linear activation function in between:

1. First Linear Transformation: embedding of the input from the model's dimensionality  $d$  to a higher-dimensional space
2. Activation Function: A non-linear function is applied to introduce non-linearity
3. Second Linear Transformation: projects the output back to the original dimensionality

$$\text{FFN}(x) = W_2 \cdot \text{ReLU}(W_1 \cdot x + b_1) + b_2$$

where  $W_1$  and  $W_2$  are weight matrices, and  $b_1$  and  $b_2$  are bias vectors.

207

## Skip connection (original paper "Deep Residual Learning for Image Recognition" 2015 by He et al.)

Skip connection

Skip connections address the vanishing (and exploding) gradient problems

- Consider a neural network layer with an input  $x$  and a desired underlying function  $H(x)$ .
- Incorporating a skip connection, the layer is restructured to model a residual  $F(x) = H(x) - x$ . Thus, the output of this layer becomes:

$$y = F(x) + x$$

### Benefits:

- Enhanced Gradient Flow
- Faster convergence and higher accuracy

208

## Multi-head attention

- Multi-Head Attention captures various aspects of the data.
- The mechanism computes multiple attention functions (heads) in parallel.
- Each head operates on linearly projected versions of the queries, keys, and values

Each head  $k$  computes the attention scores

$$\text{Attention}(Q_k, K_k, V_k) = \text{softmax} \left( \frac{XQ_k^T K_k X}{\sqrt{d_k}} \right) X V_k^T$$

where  $d_k$  is the dimension of the key vectors

### Concatenation

The outputs of all heads are concatenated and projected through a final linear layer:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)P$$

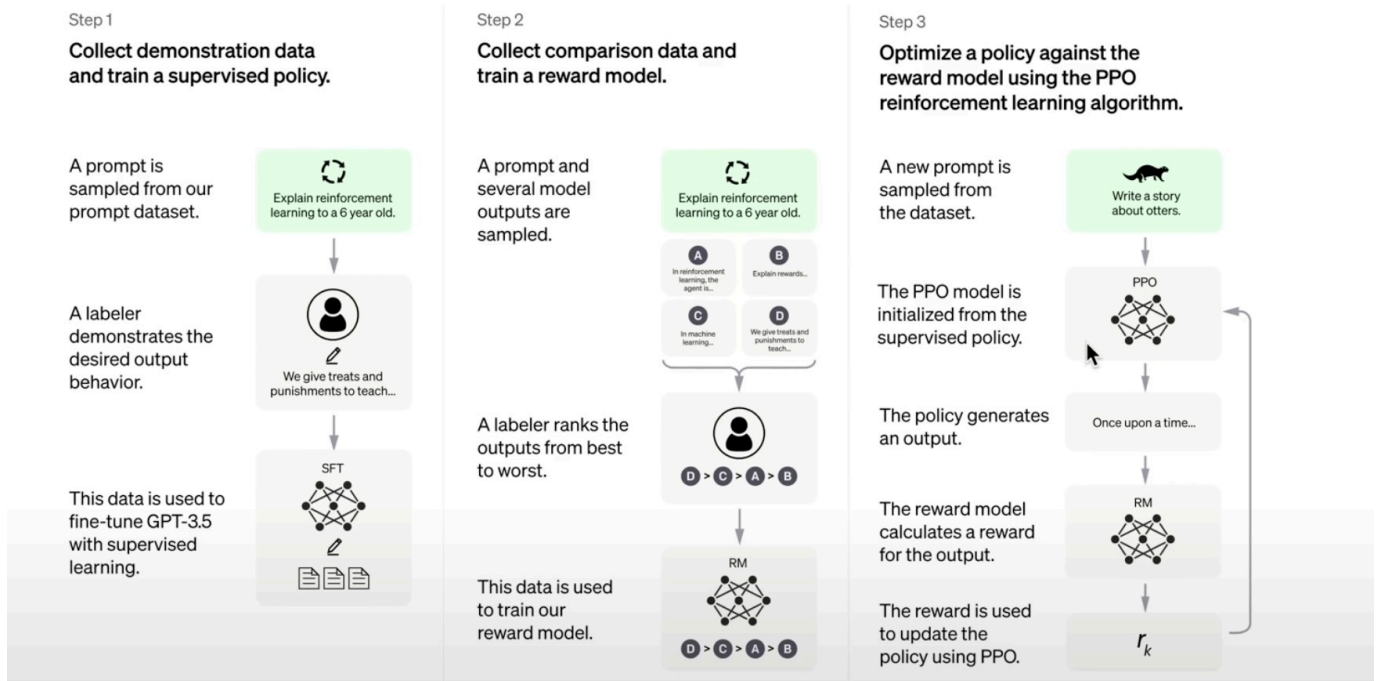
The multi head matrix is  $d \times d$ , where -  $h$  is the number of head, -  $d_k = d/h$  the size of the projection space of head k -  $P$  is a projection matrix combining the heads.

### Benefits

- Diverse Representations: Each head can capture different features or relationships within the data, enabling the model to understand various aspects of the input (grammar, style, "meaning" ...).
  - Parallel Processing: Multiple heads allow the model to process different parts of the sequence simultaneously, improving efficiency.
  - Enhanced Capacity: By attending to information from different subspaces, the model can capture complex patterns and dependencies.

# Fine tuning

After a pretraining step that resembles closely to a decoder...



ChatGPT fine tuning (from A. Karpathy lecture)

# Exercices

## Schur Complement Lemma

Let  $A$  be a square matrix partitioned as follows:

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

Assuming that  $A_{11}$  is invertible, the Schur complement of  $A_{11}$  in  $A$  is defined as:

$$S = A_{22} - A_{21}A_{11}^{-1}A_{12}$$

The lemma states that if  $A$  is invertible, then  $A$  is invertible if and only if  $S$  is invertible. Additionally, the inverse of  $A$  can be expressed as:

$$A^{-1} = \begin{bmatrix} A_{11}^{-1} + A_{11}^{-1}A_{12}S^{-1}A_{21}A_{11}^{-1} & A_{11}^{-1}A_{12}S^{-1} \\ -S^{-1}A_{21}A_{11}^{-1} & S^{-1} \end{bmatrix}$$

To demonstrate the Schur complement lemma, we can follow these



1. Start with the matrix equation  $\mathbf{AX} = \mathbf{B}$ , where  $\mathbf{A}$  is invertible and partitioned as described above.
2. Write the equation using the partitioned form of matrix  $\mathbf{A}$ :

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \end{bmatrix}$$

3. Apply block matrix operations to rewrite the equation:

$$\begin{aligned} \mathbf{A}_{11}\mathbf{X}_1 + \mathbf{A}_{12}\mathbf{X}_2 &= \mathbf{B}_1 \\ \mathbf{A}_{21}\mathbf{X}_1 + \mathbf{A}_{22}\mathbf{X}_2 &= \mathbf{B}_2 \end{aligned}$$

4. Solve the first equation for  $\mathbf{X}_1$ :

$$\mathbf{X}_1 = \mathbf{A}_{11}^{-1}(\mathbf{B}_1 - \mathbf{A}_{12}\mathbf{X}_2)$$

5. Substitute this value of  $\mathbf{X}_1$  into the second equation:

$$\mathbf{A}_{21}\mathbf{A}_{11}^{-1}(\mathbf{B}_1 - \mathbf{A}_{12}\mathbf{X}_2) + \mathbf{A}_{22}\mathbf{X}_2 = \mathbf{B}_2$$

6. Simplify the equation:

$$\mathbf{S}\mathbf{X}_2 = \mathbf{B}_2 - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{B}_1$$

7. We can see that the coefficient matrix for  $\mathbf{X}_2$  is the Schur complement,

$$\mathbf{S} = \mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12}$$

8. Therefore,  $\mathbf{X}_2$  can be calculated as:

$$\mathbf{X}_2 = \mathbf{S}^{-1}(\mathbf{B}_2 - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{B}_1)$$

9. Finally, we can substitute the values of  $\mathbf{X}_1$  and  $\mathbf{X}_2$  to obtain the solution vector  $\mathbf{X}$ :

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{11}^{-1}(\mathbf{B}_1 - \mathbf{A}_{12}\mathbf{S}^{-1}(\mathbf{B}_2 - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{B}_1)) \\ \mathbf{S}^{-1}(\mathbf{B}_2 - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{B}_1) \end{bmatrix} =$$

$$\begin{bmatrix} \mathbf{A}_{11}^{-1} & -\mathbf{A}_{11}^{-1}\mathbf{A}_{12}\mathbf{S}^{-1} & \mathbf{A}_{11}^{-1}\mathbf{A}_{12}\mathbf{S}^{-1}\mathbf{A}_{21}\mathbf{A}_{11}^{-1} & \mathbf{A}_{11}^{-1} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{S}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{0} \\ \mathbf{0} \\ \mathbf{B}_2 \end{bmatrix}$$

# Conditional Gaussian

Demonstrate the form of the distribution of  $\mathbf{x}_1 | \mathbf{x}_2$  when both vectors are gaussian.

We have the joint gaussian distribution of  $\begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix}$ :

$$\begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

We can proceed in three steps:

1. Compute  $\boldsymbol{\Sigma}^{-1} = \boldsymbol{\Lambda}$  the concentration matrix
2. Write the joint distribution into a sum of marginal in  $\mathbf{x}_2$  and conditional in  $\mathbf{x}_1$
3. Derive the conditional distribution.

## Concentration matrix

The joint covariance matrix is partitioned as:

$$\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}.$$

The concentration matrix  $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$  can be expressed as:

$$\boldsymbol{\Lambda} = \begin{bmatrix} \boldsymbol{\Lambda}_{11} & \boldsymbol{\Lambda}_{12} \\ \boldsymbol{\Lambda}_{21} & \boldsymbol{\Lambda}_{22} \end{bmatrix},$$

where the blocks of the concentration matrix are given by:

$$\begin{aligned} \boldsymbol{\Lambda}_{11} &= (\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21})^{-1}, & \boldsymbol{\Lambda}_{12} &= -\boldsymbol{\Lambda}_{11}\boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}, \\ \boldsymbol{\Lambda}_{21} &= -\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}\boldsymbol{\Lambda}_{11}, & \boldsymbol{\Lambda}_{22} &= \boldsymbol{\Sigma}_{22}^{-1} + \boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}\boldsymbol{\Lambda}_{11}\boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}. \end{aligned}$$

Here,  $\boldsymbol{\Lambda}_{11}$  is obtained using the Schur complement of  $\boldsymbol{\Sigma}_{11}$  in

$\boldsymbol{\Sigma}$ :

$$\boldsymbol{\Lambda}_{11} = (\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21})^{-1}.$$

## Decompose the joint distribution

The quadratic form of the joint distribution decomposes as follow

$$Q = (\mathbf{x}_1 - \boldsymbol{\mu}_1)^T \boldsymbol{\Lambda}_{11} (\mathbf{x}_1 - \boldsymbol{\mu}_1) + 2(\mathbf{x}_1 - \boldsymbol{\mu}_1)^T \boldsymbol{\Lambda}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2) + (\mathbf{x}_2 - \boldsymbol{\mu}_2)^T \boldsymbol{\Lambda}_{22} (\mathbf{x}_2 - \boldsymbol{\mu}_2)$$

Replacing the  $\boldsymbol{\Lambda}_{ij}$  with their expression in function of the  $\boldsymbol{\Sigma}_{ij}$ , we eventually get

$$Q = (\mathbf{x}_1 - (\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2)))^T \boldsymbol{\Sigma}_{1|2}^{-1} (\mathbf{x}_1 - (\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2))) + (\mathbf{x}_2 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2)$$

where

$$\boldsymbol{\Sigma}_{1|2} = (\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21})^{-1} = \boldsymbol{\Lambda}_{11}$$

### Derive the conditional distribution

The joint distribution of  $(\mathbf{x}_1, \mathbf{x}_2) \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , where

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix},$$

can be expressed as a product of the marginal distribution of  $\mathbf{x}_2$  and the conditional distribution of  $\mathbf{x}_1$  given  $\mathbf{x}_2$ .

The joint log-probability is:

$$\log p(\mathbf{x}_1, \mathbf{x}_2) = \log p(\mathbf{x}_2) + \log p(\mathbf{x}_1 | \mathbf{x}_2).$$

Expanding each term: 1. The marginal distribution of  $\mathbf{x}_2$ :

$$p(\mathbf{x}_2) \sim N(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22}),$$

with the log-density:

$$\log p(\mathbf{x}_2) = -\frac{1}{2} (\mathbf{x}_2 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2) - \frac{1}{2} \log |\boldsymbol{\Sigma}_{22}| - \frac{d_2}{2} \log(2\pi).$$

2. The conditional distribution of  $\mathbf{x}_1$  given  $\mathbf{x}_2$ :

$$p(\mathbf{x}_1 | \mathbf{x}_2) \sim N(\boldsymbol{\mu}_{1|2}, \boldsymbol{\Sigma}_{1|2}),$$

where

$$\boldsymbol{\mu}_{1|2} = \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2), \quad \boldsymbol{\Sigma}_{1|2} = \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21}.$$

The log-density is:

$$\log p(x_1|x_2) = -\frac{1}{2}(x_1 - \mu_{1|2})^T \Sigma_{1|2}^{-1}(x_1 - \mu_{1|2}) - \frac{1}{2} \log |\Sigma_{1|2}| - \frac{d_1}{2} \log(2\pi).$$

Thus, the joint distribution can be written as:

$$\log p(x_1, x_2) = \log p(x_2) + \log p(x_1|x_2).$$

## Stochastic Gradient and linear regression

The Mean Squared Error can be expressed as a sum of the available sample  $(y_i, \mathbf{x}_i)_i$ :

$$Q(\mathbf{w}) = \frac{1}{n} \sum_i \|y_i - \mathbf{w}^t \mathbf{x}_i\|^2 = \frac{1}{n} \sum_i Q_i(\mathbf{w})$$

### The gradient

is expressed as

$$\nabla Q_i(\mathbf{w}) = -2\mathbf{x}_i(y_i - \mathbf{w}^t \mathbf{x}_i)$$

### Online regression code

```
1 library(ggplot2)
2 library(patchwork) # Pour afficher les graphes côte à côte
3
4 # Fonction de régression linéaire séquentielle
5 regression.online <- fonction(X, Y, max.iteration = 300) {
6   p <- ncol(X)
7   n <- nrow(X)
8   X <- cbind(1, X) # Ajout de l'intercept
```

```

9
10 shuffling <- sample(1:n, n)
11 X <- X[shuffling, ]
12 Y <- Y[shuffling]
13 Q<-rep(0,max.iteration)
14 W <- rep(0, p + 1) # Initialisation des coefficients
15 for (i in 1:max.iteration) {
16   idx <- (i - 1) %% n + 1
17   x <- drop(X[idx,]) # Sélection d'un point (matrice 1x(p+1))
18   y <- Y[idx]
19   W <- W + (1 / i) * x * (y - drop(rbind(W) %*% cbind(x))) # Mise à jour du vecteur de poids
20   Q[i]<-mean((Y - X %*% cbind(W))^2)
21 }
22
23 return(list(W=W,Q=Q))
24 }

```

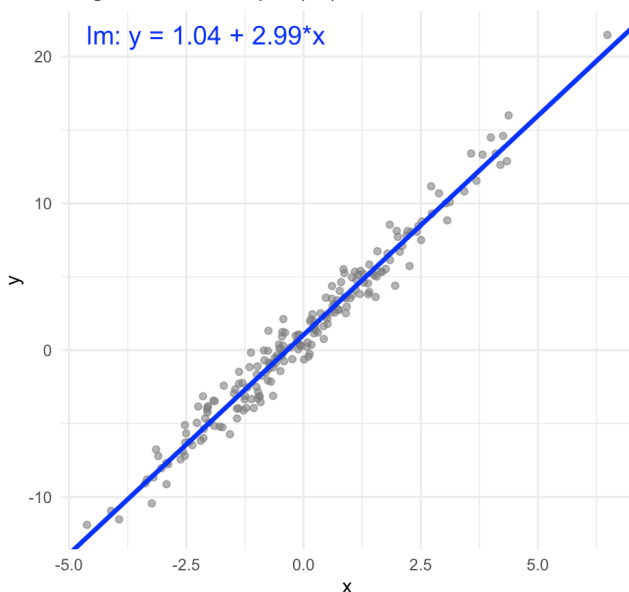
## Comparison with batch approach

```

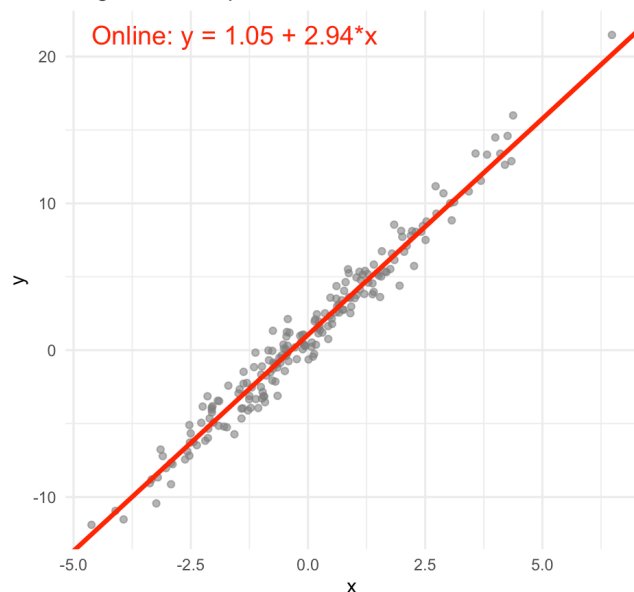
1 # Génération des données
2 set.seed(123)
3 x <- rnorm(200, sd = 2)
4 y <- 1 + 3 * x + rnorm(200, sd = 1)
5
6 # Ajustement avec lm
7 lm_model <- lm(y ~ x)
8 lm_coef <- coef(lm_model) # Coefficients de lm
9
10 # Ajustement avec la régression séquentielle
11 online_res <- regression.online(matrix(x, ncol = 1), y)
12 online_coef<-online_res$W
13
14 # Création d'un dataframe pour ggplot
15 df <- data.frame(x = x, y = y)
16
17 # Graphique 1 : Régression classique avec lm
18 p1 <- ggplot(df, aes(x = x, y = y)) +
19   geom_point(color = "gray50", alpha = 0.6) +
20   geom_abline(intercept = lm_coef[1], slope = lm_coef[2], color = "blue", lwd = 1.2) +
21   annotate("text", x = min(x), y = max(y),
22          label = sprintf("lm: y = %.2f + %.2f*x", lm_coef[1], lm_coef[2]),
23          hjust = 0, color = "blue", size = 5) +
24   labs(title = "Régression classique (lm)",
25        x = "x", y = "y") +

```

Régression classique (lm)



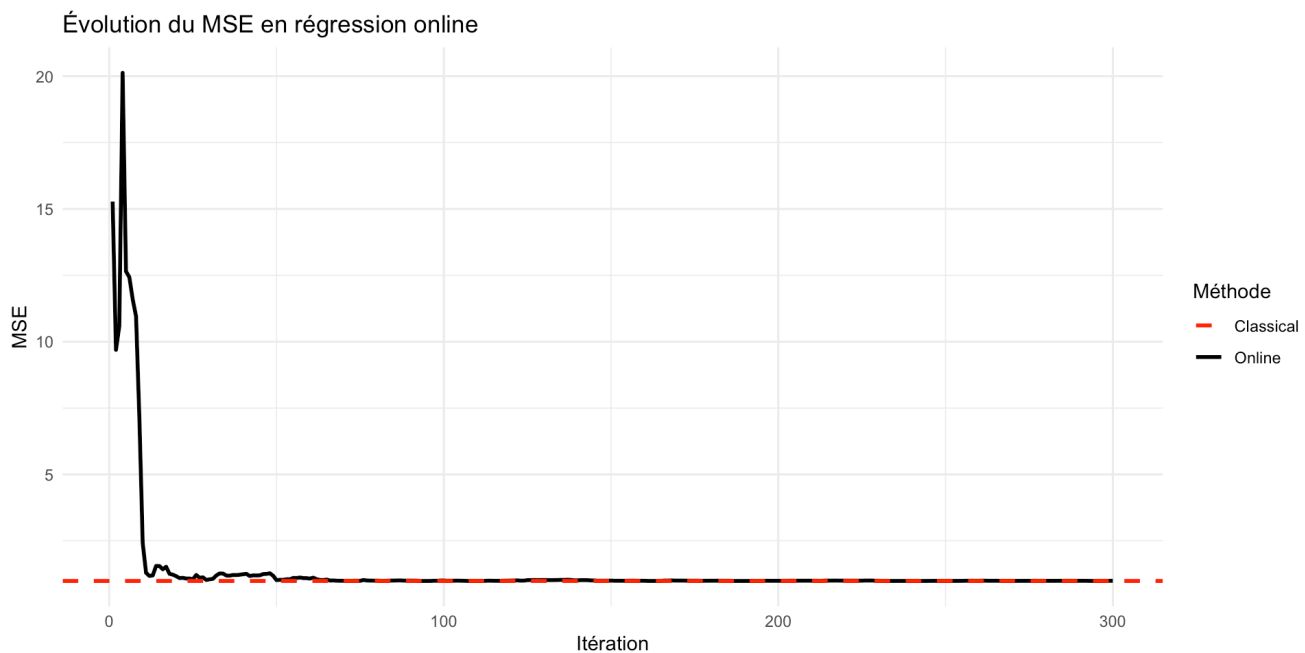
Régression séquentielle



```

1 library(ggplot2)
2
3 # Données pour la courbe de la régression online
4 df_online <- data.frame(
5   MSE = online_res$Q,
6   Iteration = 1:length(online_res$Q)
7 )
8
9 # MSE de la régression classique
10 mse_lm <- mean(lm_model$residuals^2)
11
12 # Graphique avec ggplot2
13 ggplot(df_online, aes(x = Iteration, y = MSE)) +
14   geom_line(color = "black", size = 1, aes(linetype = "Online")) + # Courbe Online en noir
15   geom_hline(aes(yintercept = mse_lm, linetype = "Classical"), color = "red", size = 1) + # Ligne lm en rouge
16   scale_linetype_manual(name = "Méthode", values = c("Online" = "solid", "Classical" = "dashed")) + # Légende
17   labs(
18     title = "Évolution du MSE en régression online",
19     x = "Itération",
20     y = "MSE"
21   ) +
22   theme_minimal()

```



# What is the KL (Kullback–Leibler) divergence between two multivariate Gaussian distributions?

KL divergence between two distributions  $P$  and  $Q$  of a continuous random variable is given by:

$$D_{KL}(p||q) = \int_x p(x) \log \frac{p(x)}{q(x)}$$

And probability density function of **multivariate Normal distribution** is given by:

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)$$

Now, let our two Normal distributions be  $\mathcal{N}(\boldsymbol{\mu}_p, \Sigma_p)$  and  $\mathcal{N}(\boldsymbol{\mu}_q, \Sigma_q)$ , both  $k$  dimensional.

$$\begin{aligned} D_{KL}(p||q) &= \mathbb{E}_p [\log(p) - \log(q)] \\ &= \mathbb{E}_p \left[ \frac{1}{2} \log \frac{|\Sigma_q|}{|\Sigma_p|} - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_p)^T \Sigma_p^{-1} (\mathbf{x} - \boldsymbol{\mu}_p) + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_q)^T \Sigma_q^{-1} (\mathbf{x} - \boldsymbol{\mu}_q) \right] \\ &= \frac{1}{2} \mathbb{E}_p \left[ \log \frac{|\Sigma_q|}{|\Sigma_p|} \right] - \frac{1}{2} \mathbb{E}_p [(\mathbf{x} - \boldsymbol{\mu}_p)^T \Sigma_p^{-1} (\mathbf{x} - \boldsymbol{\mu}_p)] + \frac{1}{2} \mathbb{E}_p [(\mathbf{x} - \boldsymbol{\mu}_q)^T \Sigma_q^{-1} (\mathbf{x} - \boldsymbol{\mu}_q)] \\ &= \frac{1}{2} \log \frac{|\Sigma_q|}{|\Sigma_p|} - \frac{1}{2} \mathbb{E}_p [(\mathbf{x} - \boldsymbol{\mu}_p)^T \Sigma_p^{-1} (\mathbf{x} - \boldsymbol{\mu}_p)] + \frac{1}{2} \mathbb{E}_p [(\mathbf{x} - \boldsymbol{\mu}_q)^T \Sigma_q^{-1} (\mathbf{x} - \boldsymbol{\mu}_q)] \end{aligned}$$

Now, since  $(\mathbf{x} - \boldsymbol{\mu}_p)^T \Sigma_p^{-1} (\mathbf{x} - \boldsymbol{\mu}_p)$  in the second term  $\in \mathbb{R}$ , we can write it as  $tr \{(\mathbf{x} - \boldsymbol{\mu}_p)^T \Sigma_p^{-1} (\mathbf{x} - \boldsymbol{\mu}_p)\}$ , where  $tr\{\}$  is the trace operator. And using the trace trick (eq 16 of section 1.1 from **Matrix Cookbook**), we can write it as  $tr \{(\mathbf{x} - \boldsymbol{\mu}_p)(\mathbf{x} - \boldsymbol{\mu}_p)^T \Sigma_p^{-1}\}$ .

The second term now is,

$$= \frac{1}{2} \mathbb{E}_p [tr \{(\mathbf{x} - \boldsymbol{\mu}_p)(\mathbf{x} - \boldsymbol{\mu}_p)^T \Sigma_p^{-1}\}]$$

The expectation and trace can be interchanged to get,

$$\begin{aligned}
&= \frac{1}{2} \text{tr} \left\{ \mathbb{E}_p \left[ (\mathbf{x} - \boldsymbol{\mu}_p)(\mathbf{x} - \boldsymbol{\mu}_p)^T \Sigma_p^{-1} \right] \right\} \\
&= \frac{1}{2} \text{tr} \left\{ \mathbb{E}_p \left[ (\mathbf{x} - \boldsymbol{\mu}_p)(\mathbf{x} - \boldsymbol{\mu}_p)^T \right] \Sigma_p^{-1} \right\}
\end{aligned}$$

We know  $\mathbb{E}_p \left[ (\mathbf{x} - \boldsymbol{\mu}_p)(\mathbf{x} - \boldsymbol{\mu}_p)^T \right] = \Sigma_p$ . Simplifying it to

$$\begin{aligned}
&= \frac{1}{2} \text{tr} \left\{ \Sigma_p \Sigma_p^{-1} \right\} \\
&= \frac{1}{2} \text{tr} \left\{ I_k \right\} \\
&= \frac{k}{2}
\end{aligned}$$

We can simplify the third term

$$\begin{aligned}
\mathbb{E}_p \left[ (\mathbf{x} - \boldsymbol{\mu}_q)^T \Sigma_q^{-1} (\mathbf{x} - \boldsymbol{\mu}_q) \right] &= \mathbb{E}_p \left[ (\mathbf{x} - \boldsymbol{\mu}_p + \boldsymbol{\mu}_p - \boldsymbol{\mu}_q)^T \Sigma_q^{-1} (\mathbf{x} - \boldsymbol{\mu}_p + \boldsymbol{\mu}_p - \boldsymbol{\mu}_q) \right] \\
&= \mathbb{E}_p \left[ (\boldsymbol{\mu}_p - \boldsymbol{\mu}_q)^T \Sigma_q^{-1} (\boldsymbol{\mu}_p - \boldsymbol{\mu}_q) + \text{tr} \left\{ \Sigma_q^{-1} \Sigma_p \right\} + 2(\boldsymbol{\mu}_p - \boldsymbol{\mu}_q)^T \Sigma_q^{-1} (\boldsymbol{\mu}_p - \boldsymbol{\mu}_q) \right] \\
&= \text{tr} \left\{ \Sigma_q^{-1} \Sigma_p \right\} + 2(\boldsymbol{\mu}_p - \boldsymbol{\mu}_q)^T \Sigma_q^{-1} (\boldsymbol{\mu}_p - \boldsymbol{\mu}_q)
\end{aligned}$$

Combining all this we get,

## Kullback–Leibler divergence between two multivariate Gaussian distributions

### Compact form

$$D_{KL}(p||q) = \frac{1}{2} \left[ \log \frac{|\Sigma_q|}{|\Sigma_p|} - k + (\boldsymbol{\mu}_p - \boldsymbol{\mu}_q)^T \Sigma_q^{-1} (\boldsymbol{\mu}_p - \boldsymbol{\mu}_q) + \text{tr} \left\{ \Sigma_q^{-1} \Sigma_p \right\} \right]$$

### When $q$ is $\mathcal{N}(0, I)$

$$D_{KL}(p||q) = \frac{1}{2} \left[ \boldsymbol{\mu}_p^T \boldsymbol{\mu}_p + \text{tr} \left\{ \Sigma_p \right\} - k - \log |\Sigma_p| \right]$$