

# Visualisation and Dimension Reduction

Christophe Ambroise

2025-02-20

## Table of contents

<b>1</b>	<b>Evaluation</b>	<b>6</b>
<b>2</b>	<b>Factor Analyser</b>	<b>7</b>
2.1	Factor Analysis . . . . .	7
2.2	The model of factor analysis . . . . .	7
2.3	Reminder: Joint and conditional Gaussian distribution (see Murphy chapter 4) . . . . .	7
2.4	Marginal and posterior distribution . . . . .	8
2.5	Estimation . . . . .	8
2.6	EM algorithm . . . . .	9
2.7	EM for factor analysis . . . . .	9
2.8	E step . . . . .	10
2.9	M step . . . . .	10
2.10	M step for $\Psi$ . . . . .	10
2.11	M step for $W$ . . . . .	11
2.12	M Step summary . . . . .	11
2.13	Implementation of the algorithm . . . . .	11
2.14	Example with the Iris . . . . .	13
2.15	Unidentifiability . . . . .	14
2.16	Possible rotations . . . . .	14
2.17	Varimax . . . . .	15
2.18	Mixture of factor analysers . . . . .	15
2.19	Relation to principal component analysis . . . . .	16
2.20	A Constrained EM Algorithm for PCA (from Ahn, J.-H. and J.-H. Oh, 2003) . . . . .	17
<b>3</b>	<b>Independent Component Analysis</b>	<b>18</b>
3.1	Independent Component analysis (Wikipedia) . . . . .	18
3.2	ICA Historical context (French Wikipedia) . . . . .	19
3.3	ICA Model . . . . .	19

3.4	Maximum likelihood estimation of ICA . . . . .	20
3.5	Estimation via Gradient Ascent . . . . .	21
3.6	Fast ICA (for one latent factor $v$ ) . . . . .	22
3.7	Fast ICA for $L > 1$ . . . . .	22
3.8	Non-Gaussianity . . . . .	23
<b>4</b>	<b>Neural networks and unsupervised learning</b>	<b>23</b>
4.1	Modeling of a neuron . . . . .	23
4.2	The neuron as a funtion . . . . .	24
4.3	Neural networks . . . . .	24
4.4	Hebbian Learning (Hebb 1949) . . . . .	25
4.5	Stochastic Gradient Descent (from Wikipedia) . . . . .	25
4.6	Iterative method . . . . .	26
4.7	Stochastic Gradient in pseudocode . . . . .	26
4.8	Adaptative learning rate . . . . .	26
4.9	K-means and Winner take all . . . . .	27
4.10	Stochastic Gradient for Kmeans (online Kmeans) . . . . .	27
4.11	Kmeans implementation . . . . .	28
4.12	One line kmeans example with Fisher iris . . . . .	28
4.13	One line kmeans example with Fisher iris . . . . .	29
4.14	PCA and Oja's rule . . . . .	29
4.15	Stochastic Gradient PCA and Oja's rule . . . . .	30
4.16	Oja's rule implementation . . . . .	31
4.17	Oja's rule example with Fisher iris . . . . .	31
<b>5</b>	<b>Variational auto-encoder</b>	<b>32</b>
5.1	Variational autoencoder ideas . . . . .	32
5.2	Auto-encoder Structure . . . . .	34
5.3	What is a VAE ? . . . . .	34
5.4	Example of use . . . . .	35
5.5	Example of use . . . . .	36
5.6	Example of use . . . . .	37
5.7	Example of use . . . . .	38
5.8	Example of use . . . . .	39
5.9	Ingredient: Parameterization of conditional distributions with Neural Networks . . . . .	39
5.10	Modeling joint distribution . . . . .	39
5.11	Modeling Conditional distribution . . . . .	39
5.12	Parameterization of conditional distributions with Neural Networks . . . . .	40
5.13	Ingredient: Stochastic Gradient . . . . .	40
5.14	What differences between Oja's rule and VAE . . . . .	40
5.15	Factor analysis generalizes Oja and is closer to a VAE . . . . .	40
5.16	Ingredient : Evidence Lower BOUND Minimization . . . . .	41
5.17	Missing data . . . . .	41

5.18	ELBO	41
5.19	Two for one	42
5.20	Alternative formulation of ELBO	42
5.21	The ELBO is maximized by Stochastic Gradient	42
5.22	Decoder Gradient: $\nabla_{\Theta} \mathcal{L}(\Theta, \Phi; x)$	43
5.23	Encoder Gradient: $\nabla_{\Phi} \mathcal{L}(\Theta, \Phi; x)$	43
5.24	Encoder function	43
5.25	Reparametrization trick	44
5.26	Computing $\log q_{\Phi}(z x)$ with a change of variable	44
5.27	Factorized Gaussian Posterior	44
5.28	Factorized Gaussian Posterior	45
5.29	Reparametrization trick	46
5.30	Full Gaussian posterior	46
5.31	Full Gaussian posterior	46
5.32	Algorithm input and output	47
5.33	Algorithm	47
5.34	Original example from Kingma: Gaussian model with MLP parametrization	48
5.35	Gaussian VAE illustrated	48
<b>6</b>	<b>Mixture of Experts</b>	<b>49</b>
6.1	Historical Context	49
6.2	MoE Components	49
6.3	Model Formulation	50
6.4	MoE Architectures	51
6.5	Example of regression mixture	51
6.6	Example of regression mixture	52
<b>7</b>	<b>Singular Value Decomposition</b>	<b>54</b>
7.1	Singular Value Decomposition	54
7.2	Applications in machine learning	54
7.3	Existence of the SVD for general matrices	55
7.4	Geometrical interpretation	55
7.5	Geometrical interpretation	56
7.6	Different versions of SVD	56
7.7	SVD $n > d$	57
7.8	SVD $n < d$	57
7.9	Existence of the SVD	57
7.10	Existence of the SVD	58
7.11	Existence of the SVD	58
7.12	Properties	59
7.13	Low rank approximation of a matrix $X$	59
7.14	A first toy example	59
7.15	Illustration of svd in image compression	61

7.16	Illustration of svd in image compression . . . . .	61
7.17	Low rank approximation of a matrix $X$ . . . . .	64
7.18	Low rank approximation of a matrix . . . . .	64
7.19	Proof . . . . .	64
7.20	Proof . . . . .	65
7.21	Proof . . . . .	65
7.22	Low rank approximation of a matrix and projection . . . . .	65
7.23	Projection on the orthogonal subspace . . . . .	66
7.24	Relation to principal component analysis . . . . .	66
7.25	Criterion . . . . .	67
7.26	Best low rank approximation . . . . .	67
7.27	Different views of the approximation . . . . .	67
7.28	Rows approximation (projection of the individuals) . . . . .	68
7.29	The approximation error . . . . .	68
7.30	Projection of the variables . . . . .	68
7.31	k first principal components . . . . .	68
7.32	Percentage of information . . . . .	69
7.33	Correlations . . . . .	69
7.34	Duality . . . . .	69
<b>8</b>	<b>Multi Dimensional Scaling</b> . . . . .	<b>69</b>
8.1	Dimensionality reduction and manifold learning . . . . .	69
8.2	A Brief history . . . . .	70
8.3	Proximity measures . . . . .	70
8.4	Classical scaling . . . . .	70
8.5	From distances to scalar product . . . . .	71
8.6	Distance and data table (continued)} . . . . .	71
8.7	Low rank approximation and principal coordinates . . . . .	72
8.8	Algorithm . . . . .	72
8.9	Optimized criterion . . . . .	73
8.10	Models and functions for the MDS . . . . .	73
8.11	In practice . . . . .	73
8.12	Error function . . . . .	73
8.13	Normalized error function . . . . .	74
8.14	Different approaches . . . . .	74
8.15	Sammon's Projection . . . . .	75
8.16	Sammon stress function . . . . .	75
8.17	Minimisation of Sammon's Stress . . . . .	75
8.18	Example in R . . . . .	76
8.19	Sammon gradient . . . . .	76
8.20	Sammon gradient descent . . . . .	76
8.21	Example Ekman colors . . . . .	77
8.22	Example Ekman colors . . . . .	77

8.23	Example Classical scaling	78
<b>9</b>	<b>t-SNE</b>	<b>78</b>
9.1	Stochastic Neighbor Embedding	79
9.2	Stochastic Neighbor Embedding	79
9.3	Cost function	80
9.4	Selecting the variance $\sigma_i$	80
9.5	Gradient descent	80
9.6	Gradient in practice	81
9.7	t-Distributed Stochastic Neighbor Embedding	81
9.8	Symmetric SNE	82
9.9	t-SNE	82
9.10	The Crowding Problem	83
9.11	Mismatched tails can compensate for mismatched dimensionalities	83
9.12	Two optimization tricks	83
9.13	Uniform manifold approximation and projection (UMAP)	84
9.14	Uniform manifold approximation and projection (UMAP)	84
9.15	UMAP vs t-SNE	84
<b>10</b>	<b>Word Embedding</b>	<b>85</b>
<b>11</b>	<b>Word embedding</b>	<b>85</b>
11.1	Words and Vectors	85
11.2	Word-document matrix	85
11.3	Word-Word matrix	85
11.4	Cosine for measuring similarity	86
11.5	TF-IDF: Weighing terms in the vector	86
11.6	TF-IDF: Weighing terms in the vector	87
11.7	Pointwise Mutual Information (PMI)	87
11.8	Positive PMI (called PPMI)	87
11.9	Applications of the tf-idf or PPMI vector models	88
11.10	Latent Semantic Analysis (LSA)	88
11.11	Latent Semantic Analysis	89
11.12	Latent Semantic Analysis	89
11.13	Latent Dirichlet Allocation (LDA)	90
11.14	LDA Generative process	90
11.15	Latent Dirichlet Allocation (LDA) in summary	91
11.16	Latent Dirichlet Allocation (LDA) estimation	91
11.17	Example of LDA in Action	91
11.18	Pro and Cons	92
11.19	Word2vec	92
11.20	Embedding derived from classification	93
11.21	From dot product to logistic regression	93

11.22	Example	94
11.23	Skip-gram model learns two separate embeddings	95
11.24	Final embeddings	95
<b>12</b>	<b>Transformers</b>	<b>95</b>
12.1	The nature of transformers	95
12.2	The nature of transformers	96
12.3	Optimized criteria	97
12.4	Masked Language Modeling (MLM)	97
12.5	Conditional Language Modeling / Translation Loss	98
12.6	Initial embedding	98
12.7	A minimal self-attention architecture for decoder	99
12.8	Attention and matrix notation	99
12.9	Illustration of self-Attention	100
12.10	Interpretation of Q, K, V	100
12.11	Important details	100
12.12	Positional encoding	101
12.13	Feed Forward Neural Network	101
12.14	Skip connection (original paper “Deep Residual Learning for Image Recognition” 2015 by He et al.)	102
12.15	Multi-head attention	103
12.16	Fine tuning	104
<b>13</b>	<b>Exercices</b>	<b>104</b>
13.1	Schur Complement Lemma	104
13.2	Conditional Gaussian	106
13.3	Stochastic Gradient and linear regression	108
13.4	What is the KL (Kullback–Leibler) divergence between two multivariate Gaussian distributions?	112
13.5	Kullback–Leibler divergence between two multivariate Gaussian distributions	114

## 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

### 2.1 Factor Analysis

- Using discrete latent variables provides limited summary (clustering)
- An alternative is to use a vector of real-valued latent variables,  $z \in \mathbb{R}^L$ .
- “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.

### 2.2 The model of factor analysis

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

$$x = Wz + \mu + \epsilon$$

where

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

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

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

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

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

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

Let us recall that if

- $z \sim \mathcal{N}(\mu_z, \Sigma_{zz})$  and  $x \sim \mathcal{N}(\mu_x, \Sigma_{xx})$  and  $cov(z, x) = \Sigma_{zx}$

then

$$p(z, x) = N\left(\begin{bmatrix} z \\ x \end{bmatrix} \mid \begin{bmatrix} \mu_z \\ \mu_x \end{bmatrix}, \begin{bmatrix} \Sigma_{zz} & \Sigma_{zx} \\ \Sigma_{xz} & \Sigma_{xx} \end{bmatrix}\right)$$

and

$$p(z, x) = p(z|x)p(x) = \mathcal{N}(z|\mu_{z|x}, \Sigma_{z|x})\mathcal{N}(x|\mu_x, \Sigma_{xx})$$

where

- $\mu_{z|x} = \mu_z + \Sigma_{zx}\Sigma_{xx}^{-1}(x - \mu_x)$
- $\Sigma_{z|x} = \Sigma_{zz} - \Sigma_{zx}\Sigma_{xx}^{-1}\Sigma_{xz}$

## 2.4 Marginal and posterior distribution

### Marginal distribution

$$x \sim \mathcal{N}_D(\mu, \Sigma_{xx} = WW^T + \Psi)$$

### Posterior distribution

$$z|x \sim \mathcal{N}_L(\mu_{z|x}, \Sigma_{z|x})$$

where

- $\Sigma_{z|x} = (I_L + W^T\Psi^{-1}W)^{-1} = S$
- $\mu_{z|x} = \Sigma_{z|x}\Sigma_{xx}^{-1}(x - \mu) = SW^T\Psi^{-1}(x - \mu)$

### Exercise

Demonstrate the above formulas

## 2.5 Estimation

### The mean $\mu$

can be estimated by maximum likelihood

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



$W$  and  $\Psi$

are estimated using an EM algorithm

## 2.6 EM algorithm

### Data

- Observed data :  $x_{1:n}$
- Missing (or hidden) data :  $z_{1:n}$

### Principle

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

## 2.7 EM for factor analysis

Let us assume that  $\mu = 0$  (centering of the  $x_i$ ), the complete log-likelihood is

$$\begin{aligned}\log p(X, Z|\mu, W, \Psi) &= \sum_i \log \mathcal{N}_L(z_i; 0, I) + \log \mathcal{N}_D(x_i; Wz_i, \Psi) \\ &= -\frac{n}{2} \log |I_L| - \frac{n}{2} \operatorname{Tr}(\hat{\Sigma}_{zz}) \\ &\quad -\frac{n}{2} \log |\Psi| - \frac{n}{2} \operatorname{Tr}(\hat{\Sigma}_{xx} \Psi^{-1}) + Cst\end{aligned}$$

where

$$\hat{\Sigma}_{xx} = \frac{1}{n} \sum_i (x_i - Wz_i)(x_i - Wz_i)^T$$

### Exercise

Demonstrate the above formula

## 2.8 E step

The expectation of the complete log-likelihood requires

1.  $\mathbb{E}_{z|x}[z_i] = SW^T\Psi^{-1}(x_i - \mu)$  where  $S = (I_L + W^T\Psi^{-1}W)^{-1}$
2.  $\mathbb{E}_{z|x}[z_i z_i^T] = E_{z|x}[z_i]E_{z|x}[z_i^T] + S$

## 2.9 M step

Reminders

$$\begin{aligned}\frac{\partial(b^T a)}{\partial a} &= b \\ \frac{\partial(a^T A a)}{\partial a} &= (A + A^T)a \\ \frac{\partial}{\partial A} \text{tr}(BA) &= B^T \\ \frac{\partial}{\partial A} \log |A| &= (A^{-1})^T \\ \text{tr}(ABC) &= \text{tr}(CAB) = \text{tr}(BCA)\end{aligned}$$

Thus if  $x$  is a vector

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

## 2.10 M step for $\Psi$

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

where

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

## 2.11 M step for $W$

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

## 2.12 M Step summary

### Loading matrix

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

### Noise covariance matrix

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

### Log-likelihood

The log-likelihood can be computed using the EM decomposition

$$\log P(X; \Theta) = E_{Z_{1:n}|x_{1:n}} [\log P(x_{1:n}, z_{1:n}; \theta)] - E_{Z_{1:n}|x_{1:n}} [\log P(z_{1:n}|x_{1:n}; \theta)]$$

## 2.13 Implementation of the algorithm

### Initialisation via a PCA

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

## E step

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

## M Step

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

## Computation of the criterion

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

## Putting it all together

```

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

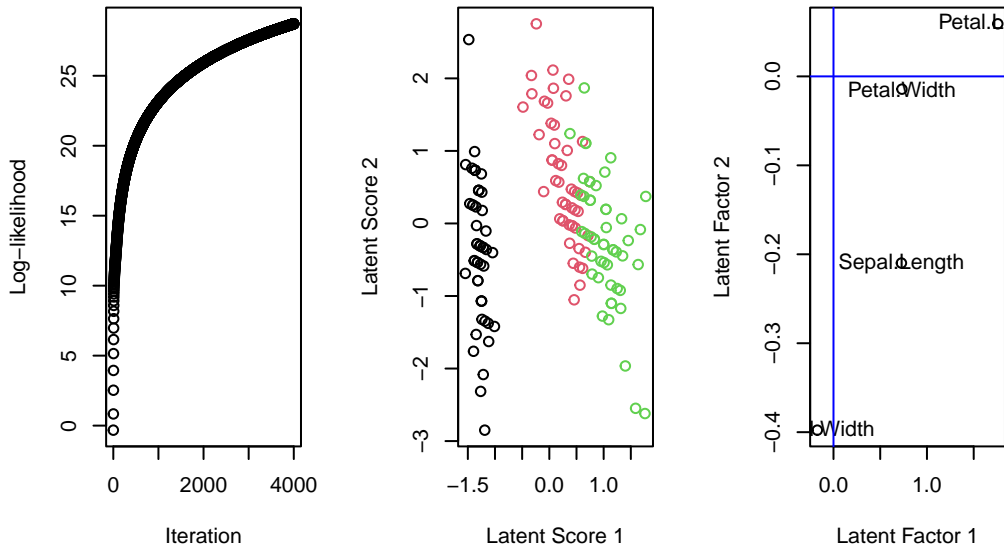
```

## 2.14 Example with the Iris

```

data(iris); L=2
X<-scale(iris[,1:4],center=TRUE,scale=FALSE)
FA.result<-FA.EM(X,L=L,max.iter=4000)
names(FA.result$W)<-paste("Latent Factor",1:L)
names(FA.result$M)<-paste("Latent Score",1:L)
par(mfrow=c(1,3))
plot(FA.result$log.likelihood,ylab="Log-likelihood",xlab="Iteration")
plot(FA.result$M,col=iris$Species)
plot(FA.result$W)
text(FA.result$W,labels = names(iris)[1:4])
abline(h=0,v=0,col="blue")

```



## 2.15 Unidentifiability

- If we consider  $R$  an orthogonal rotation matrix such that

$$RR^T = I$$

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

- $W$  cannot be uniquely identified.

## 2.16 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.

## 2.17 Varimax

```
W.FA<-FA.result$W
W.Varimax<-varimax(as.matrix(FA.result$W))$loadings
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

## 2.18 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:

$$p(x_i | z_i, q_i = k) = \mathcal{N}(x_i | \mu_k + W_k z_i, \Psi)$$

where  $k$  is the component number and  $W_k$  is a loading matrix defining the relation between the observation  $x_i$  and the latent vector  $z_i$

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

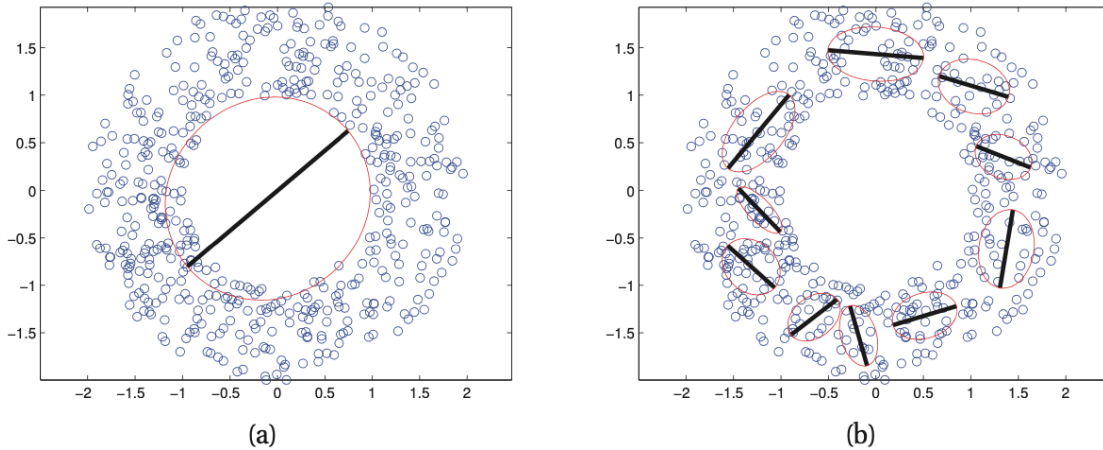


Figure 1: Mixture of 1d PPCAs with 1 and 10 components (from Murphy Chapter 12)

## 2.19 Relation to principal component analysis

### Assumption

If

- $\Psi = \sigma^2 I$
- $W$  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

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

where  $W^T W = I$



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

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

```

```
library(ggplot2)
library(cowplot)
#library(ggfortify)
library(metR)
X<-scale(iris[,1:4])
W <- data.frame(prcomp(X, scale. = TRUE)$rotation)

iris.pca<-ggplot(data=W,aes(PC1,PC2,label=rownames(W)))+geom_text(col="red")+
  geom_vline(aes(xintercept=0))+
  geom_hline(aes(yintercept=0))+
  geom_segment(aes(x=0, y=0, xend=PC1, yend=PC2 ), arrow=arrow(length=unit(0.3,"cm")), col="red")

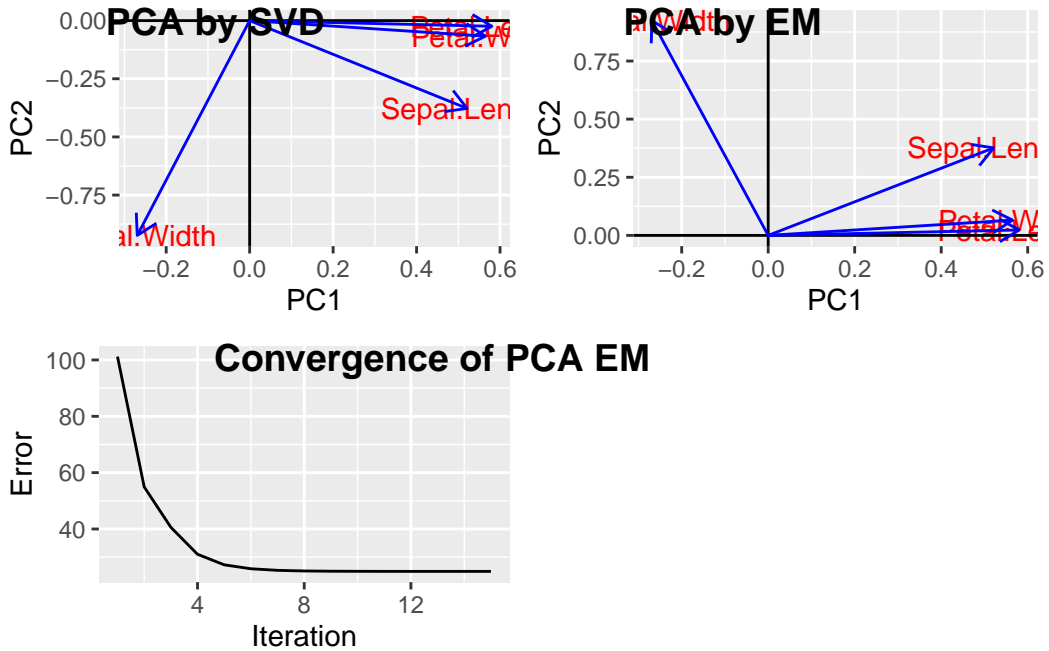
pca.em.res<-PCA.EM(X)
W<-pca.em.res$W
names(W)<-paste("PC",1:2,sep="")

iris.em.pca<-ggplot(data=W,aes(PC1,PC2,label=rownames(W)))+geom_text(col="red")+
  geom_vline(aes(xintercept=0))+
  geom_hline(aes(yintercept=0))+
  geom_segment(aes(x=0, y=0, xend=PC1, yend=PC2 ), arrow=arrow(length=unit(0.3,"cm")), col="red")

Error=data.frame(Error=pca.em.res$Error,Iteration=1:length(pca.em.res$Error))
iris.em.pca.error<-ggplot(data=Error, aes(x=Iteration,y=Error))+geom_line()

```

```
plot_grid(iris.pca,iris.em.pca,
          iris.em.pca.error,
          labels = c("PCA by SVD","PCA by EM","Convergence of PCA EM"))
```



### 3 Independent Component Analysis

#### 3.1 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.

## 3.2 ICA Historical context (French Wikipedia)

### Blind source separation

The first formulation was carried out in 1984 by Jeanny Héroult 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

- A first formalism of the blind source separation problem, as well as an algorithm making it possible to obtain a solution, was proposed by C. Jutten and J. Héroult in 1991.
- Mathematical formalization in the simplest case (linear mixing snapshot) was carried out in 1994 by P. Comon

## 3.3 ICA Model

Let  $x_t \in \mathbb{R}^D$  be the observed signal at the sensors at ‘time’  $t$ , and  $z_t \in \mathbb{R}^L$  be the vector of source signals:

$$x_t = Wz_t + \epsilon_t$$

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

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

However, we will use a **different prior** for  $p(z_t)$ .

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

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

$$p(z_t) = \prod_j p_j(z_{tj}).$$

### Additional assumptions

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

## 3.4 Maximum likelihood estimation of ICA

If the data is centered and whitened, we have

$$\mathbb{E}[xx^t] = I = W\mathbb{E}[zz^T]W^T = WW^T$$

also have

Hence we see that  $W$  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  $V = W^{-1}$  these are often called the recognition weights, as opposed to  $W$ , which are the generative weights

## Log-likelihood

Since  $x = Wz$ , we have

$$p_x(Wz_t) = p_z(z_t)|\det(W^{-1})| = p_z(Vx_t)|\det(V)|$$

Hence assuming T iid samples:

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

since  $V$  is orthogonal.

## 3.5 Estimation via Gradient Ascent

Let us define  $h_j = v_j^T x$ ,

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

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

Repeat for each datapoint  $x$ :

1. Put  $x$  through a linear mapping:

$$h = Vx$$

2. Put  $h$  through a nonlinear map:

$$g_j = g_j(h_j)$$

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

3. Adjust the weights in accordance with

$$\nabla V \propto [V^T]^{-1} + xg^T.$$

- matrix inversion results in a slow algorithm

### 3.6 Fast ICA (for one latent factor $v$ )

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

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

Let us make the approximation

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

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

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

Which can be expressed as

$$v := E[xg(v^T x)] - E[g'(v^T x)]v$$

(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

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

### 3.7 Fast ICA for $L > 1$

**Centering and Withining is assumed**

For  $L > 1$  the process is iterated for all  $v_j$  with orthogonalisation

## Fast ICA for $L > 1$

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

Notice that the expectation can be approximated with monte-carlo (choosing for example 1 data-point...)

## 3.8 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  $g(u) = u^3$

## 4 Neural networks and unsupervised learning

### 4.1 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

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

## Dendrites

are characterized by a weight vector  $w$

## The output

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

## 4.2 The neuron as a function

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

$$h(x, w) = x^T w$$

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

A typical activation function is for example:

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

The functions which have this appearance are said to be sigmoid

## 4.3 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)



## 4.4 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.

## 4.5 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:

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

where the parameter  $w$  that minimizes  $Q(w)$  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:

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

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

## 4.6 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:

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

- 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.

## 4.7 Stochastic Gradient in pseudocode

1. Choose an initial vector of parameters  $w$  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:  $w := w - \eta \nabla Q_i(w)$ .

## 4.8 Adaptive 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.

## 4.9 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,  $w_{i^*}$ , checks:

$$\forall i, (w_{i^*} \cdot x) \geq (w_i \cdot x).$$

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

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

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

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

## 4.10 Stochastic Gradient for Kmeans (online Kmeans)

The criterion to be optimized can be written as

$$Q(w_1, \dots, w_K) = \frac{1}{2n} \sum_i \sum_k I_{(k=\arg \min_{\ell} \|x_i - w_{\ell}\|^2)} \|x_i - w_k\|^2$$

1. Choose an initial vector of parameters  $w$  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  $k = 1, 2, \dots, K$ 
      - $\nabla Q_i(w_k) = -\frac{1}{n} I_{(k=\arg \min_{\ell} \|x_i - w_{\ell}\|^2)} (x_i - w_k)$
      - $w_k := w_k - \eta \nabla Q_i(w_k)$ .

## 4.11 Kmeans implementation

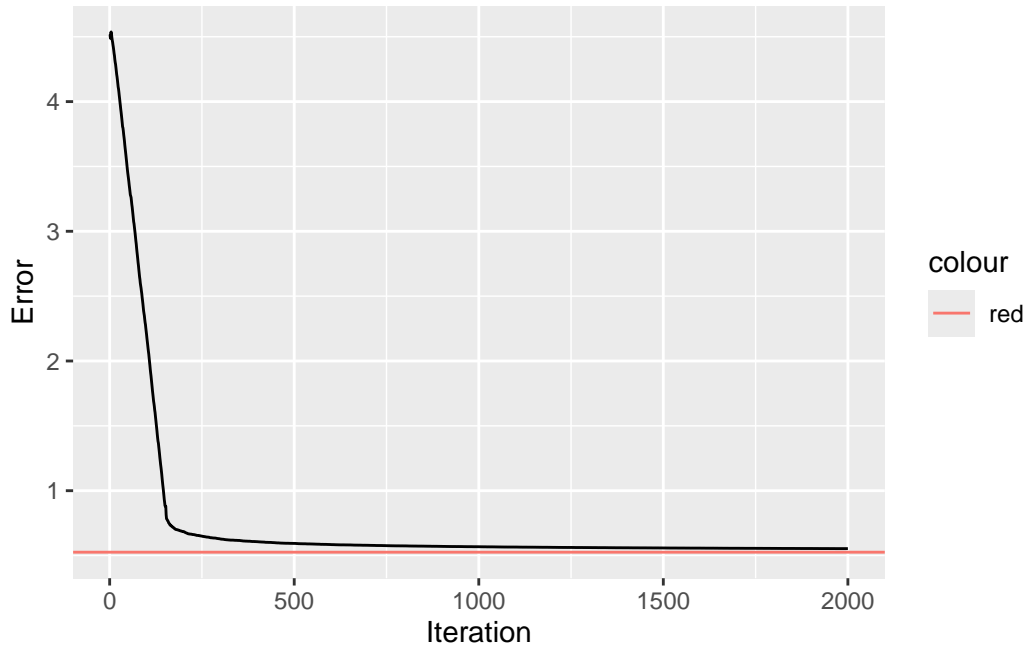
```
kmeans.winner.take.all<-function(X,K=2,max.iteration=2000){
  p<-ncol(X); n<-nrow(X);shuffling<-sample(1:n,n)
  X<-X[shuffling,]; W<-X[sample(1:n,K),]
  Q<-rep(0,max.iteration); cluster<-rep(0,n)
  distances<-rep(sum(diag(var(X)))*(n-1)/n,n)
  for (i in 1:max.iteration){
    x<-cbind(X[(i-1)%n +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

## 4.12 One line kmeans example with Fisher iris

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

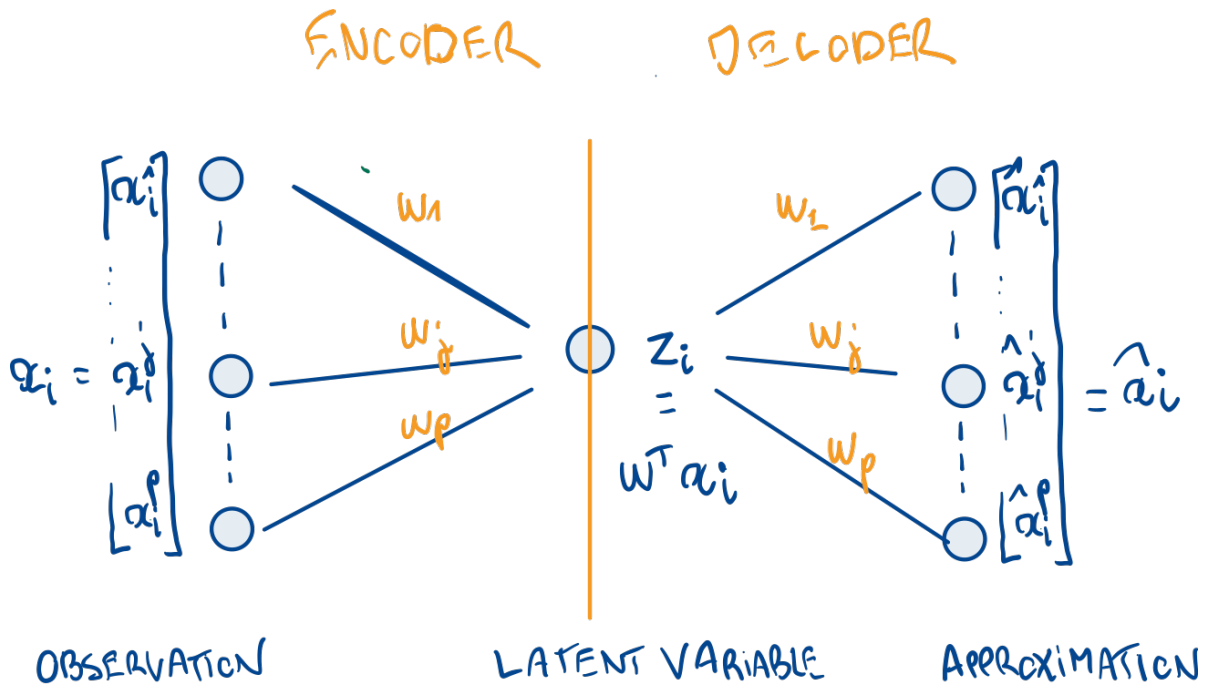
	setosa	versicolor	virginica
1	50	0	0
2	0	45	11
3	0	5	39

### 4.13 One line kmeans example with Fisher iris



### 4.14 PCA and Oja's rule

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



$$\hat{a}_i = W^T a_i W$$

OJA ESTIMATES  $W$

Figure 2: PCA according OJA

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

$$w := w - \eta z(x - zw)$$

#### 4.15 Stochastic Gradient PCA and Oja's rule

The criterion to be optimized can be written as

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

where  $\|w\|^2 = 1$ .

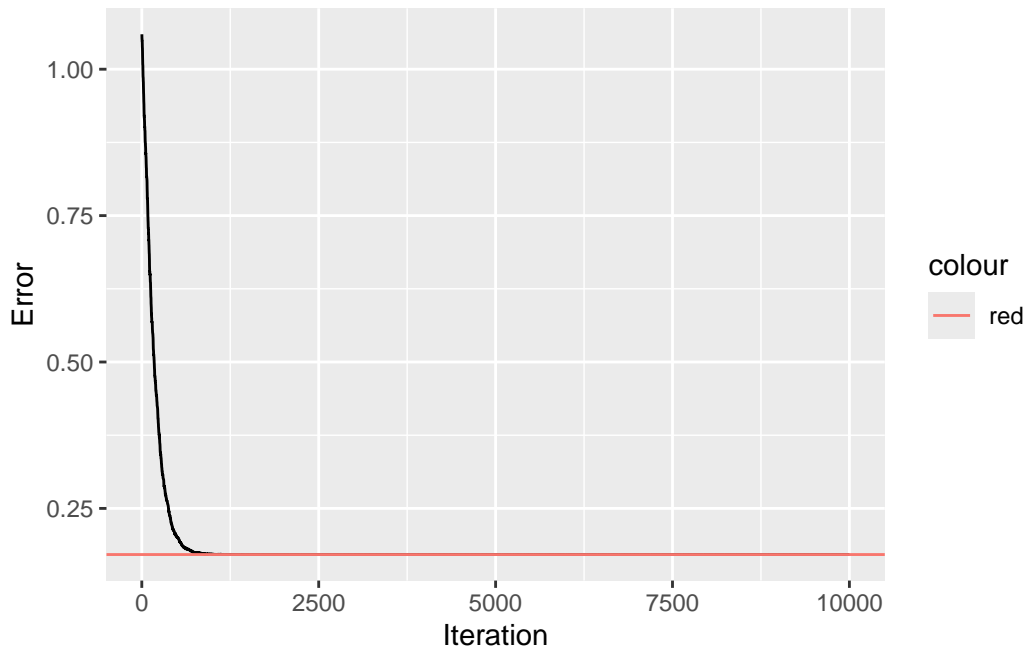
1. Choose an initial vector of parameters  $w$  and learning rate  $\eta$
2. Repeat until an approximate minimum is obtained:
  - a. Randomly shuffle examples in the training set.
  - b. For  $i = 1, 2, \dots, n$  do:
    - $\nabla Q_i(w) = y_i(x_i - y_i w)$
    - $w := w - \eta \nabla Q_i(w)$ .

## 4.16 Oja's rule implementation

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

## 4.17 Oja's rule example with Fisher iris

```
data(iris)
X<-scale(iris[,1:4],center=TRUE,scale=FALSE)
n<-nrow(X)
X<-as.matrix(X[sample(1:n),])
pca.oja<-Oja.rule(X)
w.oja<-pca.oja$w
Q<-pca.oja$Q
w.pca<-rbind(eigen(cov(X))$vector[,1])
Error.pca<- 1/(2*n) *sum((X - X%*%t(w.pca)%*%w.pca)^2)
library(ggplot2)
ggplot(data=data.frame(Error=Q,Iteration=1:length(Q)),
       aes(x=Iteration,y=Error))+geom_line()+
  geom_hline(aes(yintercept=Error.pca,col="red"))
```



## 5 Variational auto-encoder

### 5.1 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 it 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

## 5.2 Auto-encoder Structure

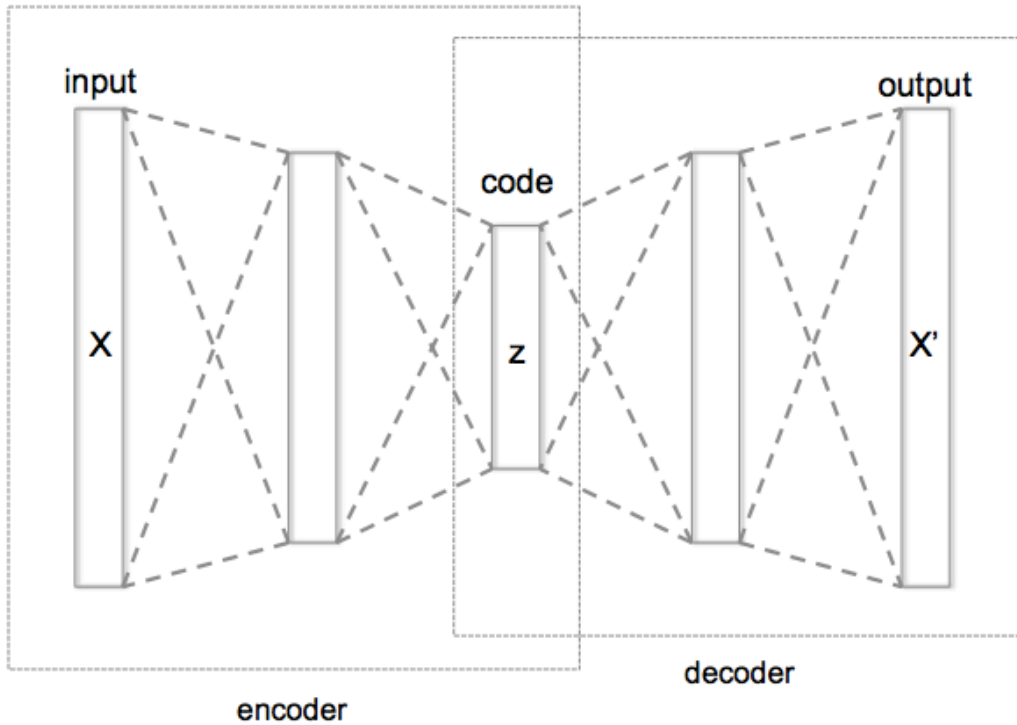


Figure 3: Auto-encoder Structure

## 5.3 What is a VAE ?

### Coupling of 2 parametric models

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

$$p(x) = \int p(x, z) dz$$

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

### 5.4 Example of use

#### Data



Left batch of original training set - right : random generation of images

## 5.5 Example of use

### Learning from data

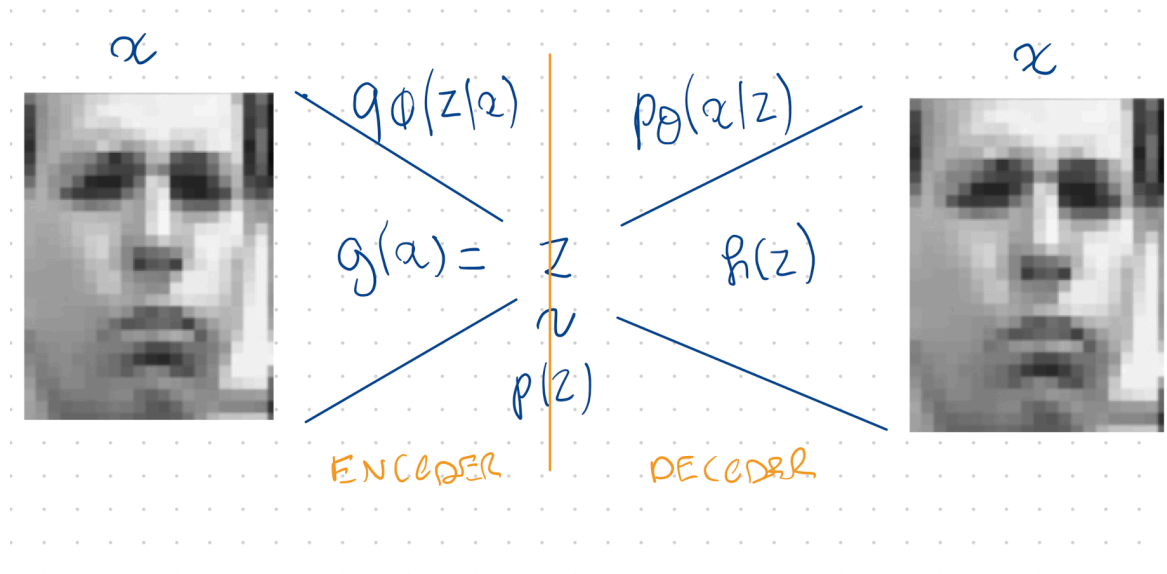


Figure 4: Frey image learning

## 5.6 Example of use

### New Data generation from latent simulation

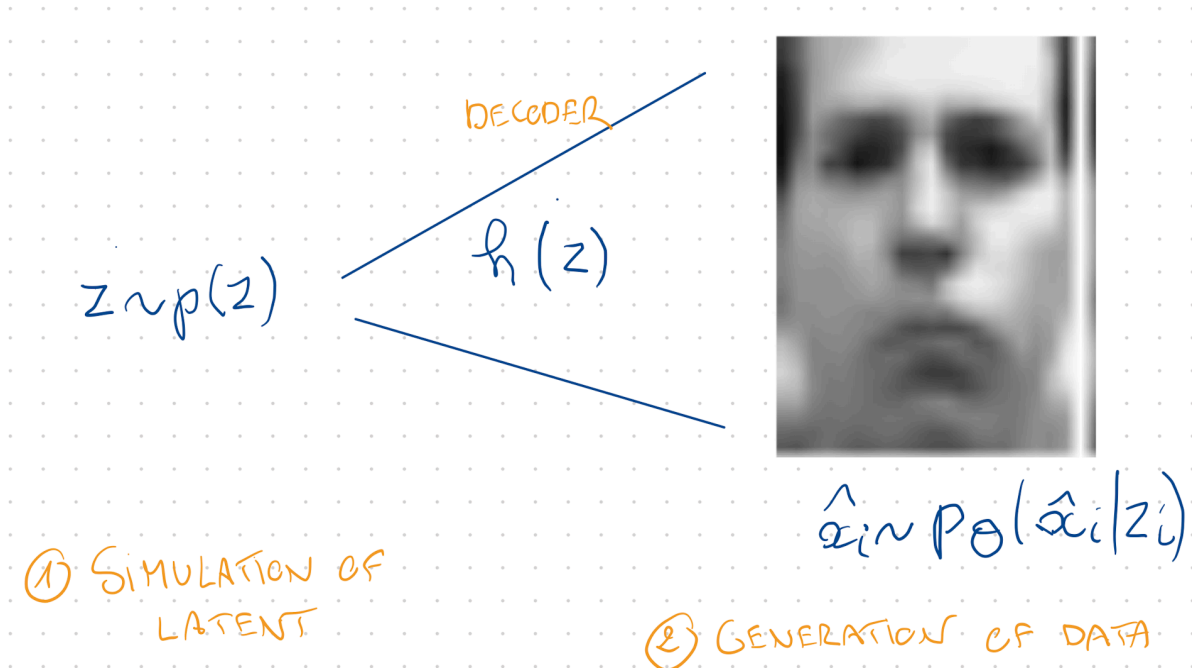


Figure 5: Frey image generation

## 5.7 Example of use

### Data representation in latent space

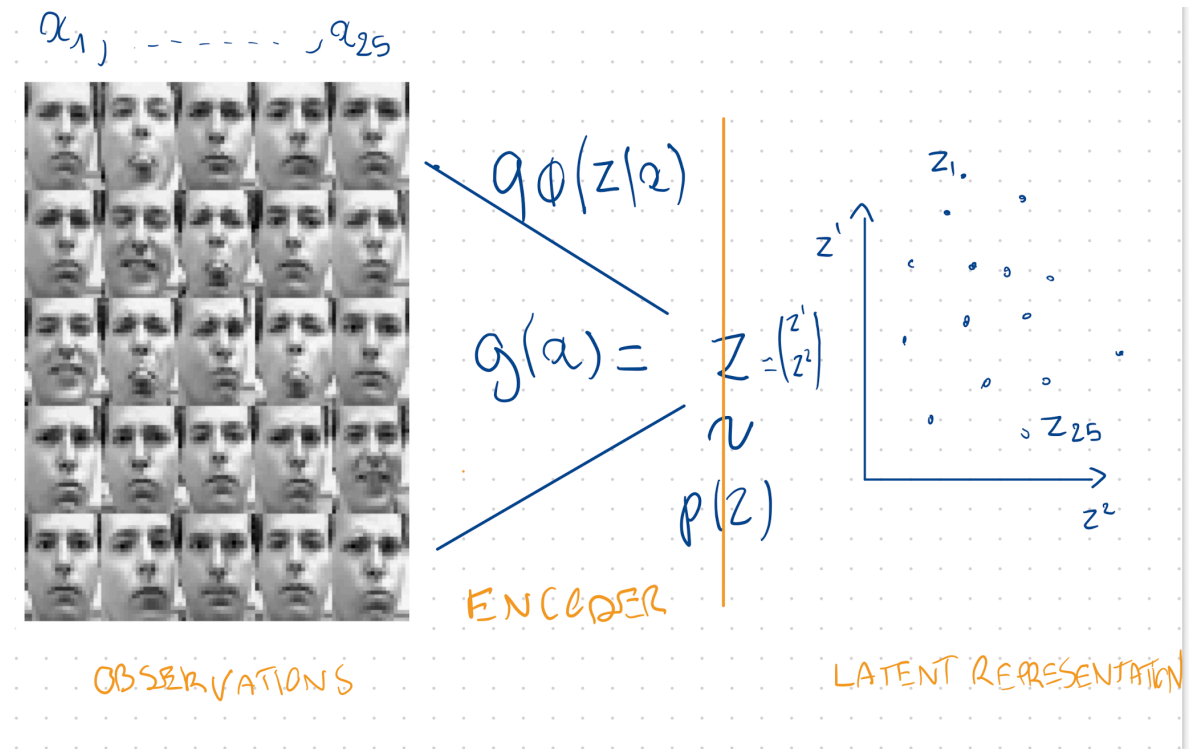


Figure 6: Frey image representation

## 5.8 Example of use

### Missing data imputation after learning

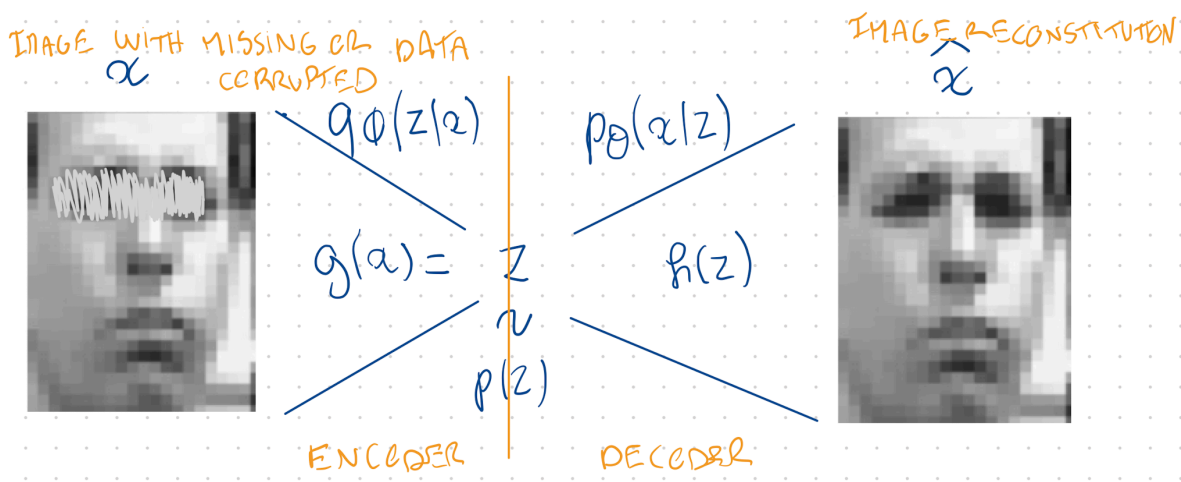


Figure 7: Frey image pixel imputation

## 5.9 Ingredient: Parameterization of conditional distributions with Neural Networks

### 5.10 Modeling joint distribution

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

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

### 5.11 Modeling Conditional distribution

#### Classification and regression

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

## 5.12 Parameterization of conditional distributions with Neural Networks

### Classification

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

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

## 5.13 Ingredient: Stochastic Gradient

## 5.14 What differences between Oja's rule and VAE

### What is a VAE ?

#### Coupling of 2 parametric models

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

### In Oja's rule

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

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

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

$$x = Wz + \mu + \epsilon$$

where

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



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

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

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

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

## 5.16 Ingredient : Evidence Lower Bound Minimization

### 5.17 Missing data

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

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

#### Approximation

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

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

### 5.18 ELBO

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

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

## 5.19 Two for one

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

## 5.20 Alternative formulation of ELBO

ELBO can be rewritten as

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

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

### Exercises

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

## 5.21 The ELBO is maximized by Stochastic Gradient

Let  $X$  be a i.i.d sample of random vector  $x$

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

### Gradient

The Gradient can be separated into 2 parts

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

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

Given some usually verified conditions and a Monte Carlo Approximation

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

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

Encoder Gradient is more difficult to compute since in general

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

A reparametrization (variable change) trick allows a workaround

## 5.24 Encoder function

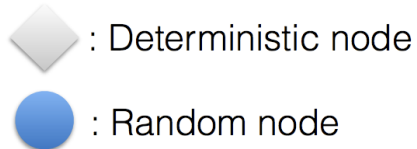
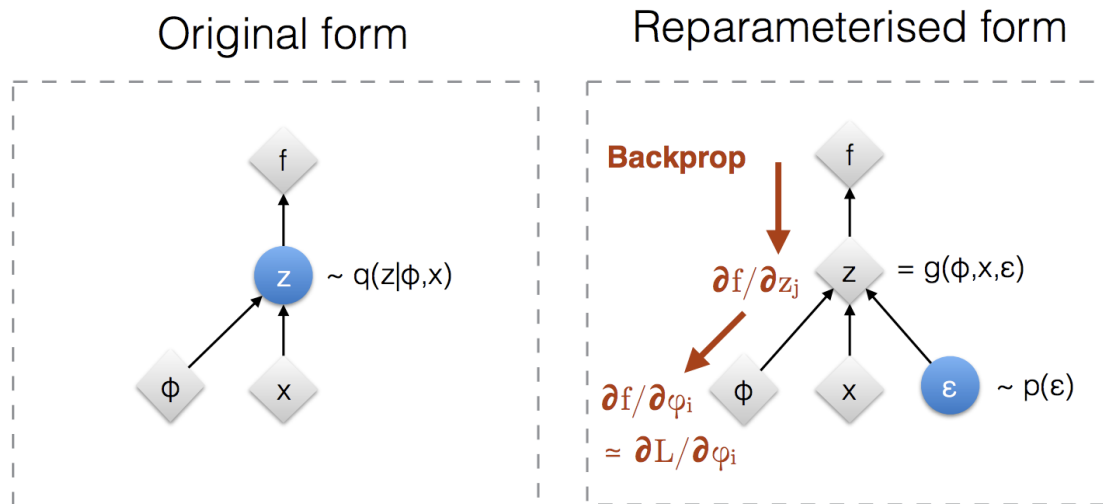
Let us rewrite the decoder function with a random vector  $\epsilon$  whose distribution is not parametrized by  $\Phi$ :

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

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

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

## 5.25 Reparametrization trick



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

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

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

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

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

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

## 5.27 Factorized Gaussian Posterior

**Model**

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

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

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

## Reparametrization

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

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

## 5.28 Factorized Gaussian Posterior

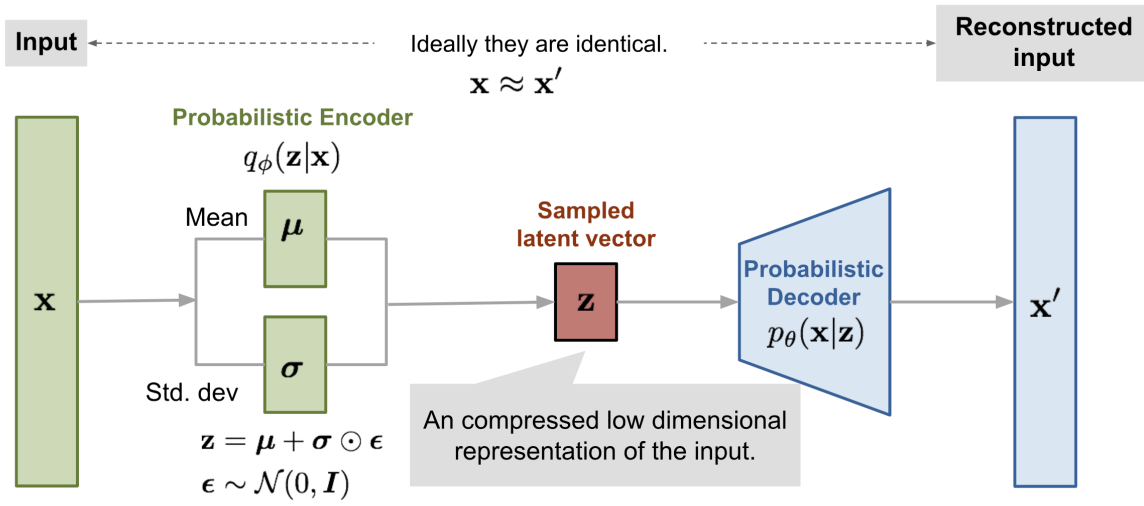
The Jacobian of the transformation is

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

The log posterior density is

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

## 5.29 Reparametrization trick



## 5.30 Full Gaussian posterior

### Model

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

### Reparametrization

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

$$z = \mu + L\epsilon$$

Where  $L$  is a lower triangular matrix obtained from a Cholesky decomposition of  $\Sigma = LL^T$

## 5.31 Full Gaussian posterior

The Jacobian has a simple form

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

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

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

### 5.32 Algorithm input and output

Input:

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

Output:

- $\Theta$
- $\Phi$

### 5.33 Algorithm

Initialisation of  $\Theta$  and  $\Phi$

While SGD not converged do

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

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

– its gradients

$$\begin{aligned}
* \nabla_{\Theta} \tilde{\mathcal{L}}_{\Theta, \Phi}(X^M, \epsilon) &= \frac{1}{m} \sum_{x \in X^M} \frac{\partial \log p_{\Theta}(x, z)}{\partial \Theta} \\
* \nabla_{\Phi} \tilde{\mathcal{L}}_{\Theta, \Phi}(X^M, \epsilon) &= \frac{1}{m} \sum_{x \in X^M} \frac{\log \partial d_{\Phi}(x, \epsilon)}{\partial \Phi}
\end{aligned}$$

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

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

### Multivariate Gaussian decoder with a diagonal covariance structure

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

#### Decoder

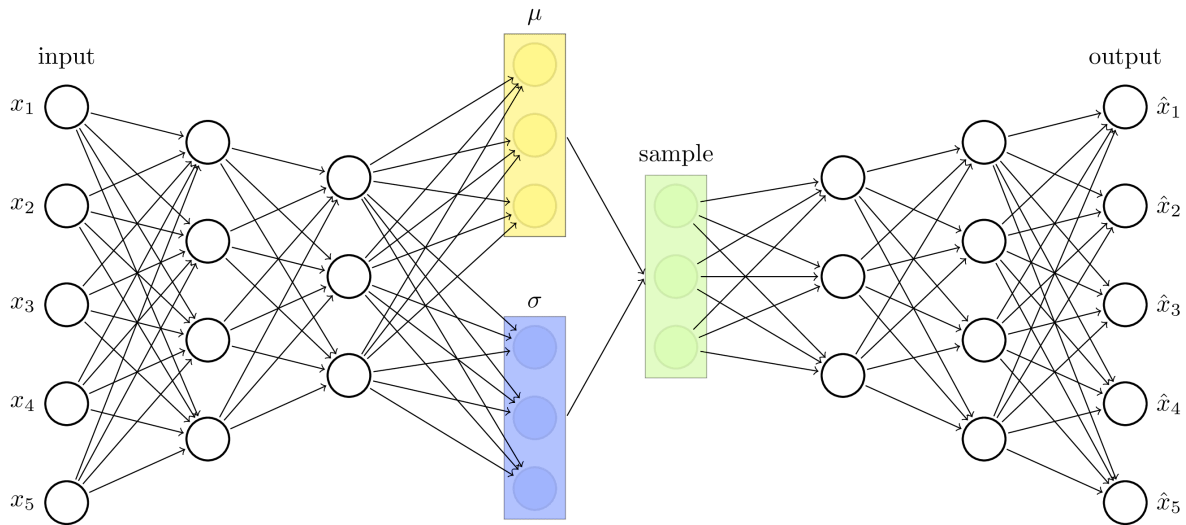
- $\log p(x|z) = \log \mathcal{N}(x; \mu, \sigma^2 I)$  where  $\mu = W_4 h + b_4$
- $\log \sigma^2 = W_5 h + b_5$
- $h = \tanh(W_3 z + b_3)$

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.

#### Encoder

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

## 5.35 Gaussian VAE illustrated





## 6 Mixture of Experts

### 6.1 Historical Context

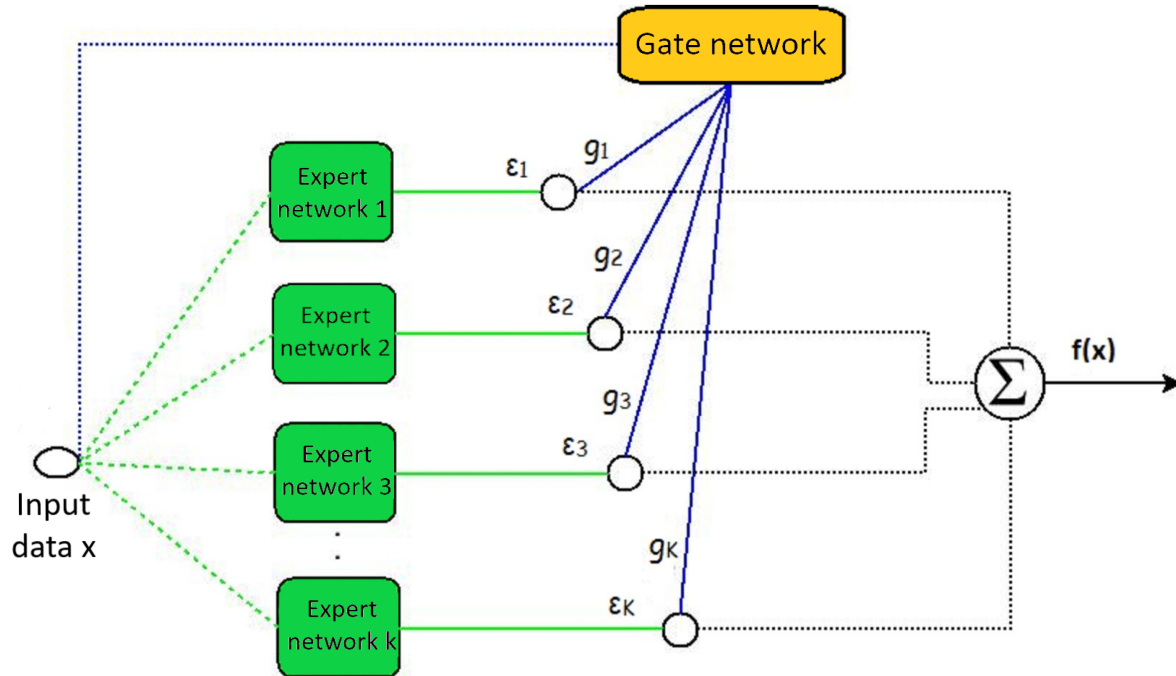
- Introduced in the early 1990s by Jacobs et al. (1991).

#### 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 sub-tasks, improving prediction accuracy (Masoudnia & Ebrahimpour, 2014; Yuksel et al., 2012).
- **Widely applied in machine learning**, including speech recognition, image processing, LLM (Deepseek, Mixtral, ChatGPT...)

### 6.2 MoE Components

- **Experts**: Independent models (e.g., neural networks, regressions) specialized in subsets of data.
- **Gating Network**: Evaluates input data and assigns relevance weights to experts.
- Combines expert outputs into a final prediction (Gormley & Eisner, 2019).



### 6.3 Model Formulation

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

- $g_k(x_i)$ : Output of the gating network

#### Probabilistic formulation

- $f(y_i; x_i, \Theta) = p(y_i | x_i, \Theta)$
- the gating network comes from a latent variable  $z$ :

$$g_k(x_i) = p(z_i = k | x_i)$$

where

$$z_i = k | x_i \sim M(1; g(x_i) = (g_1(x_i), \dots, g_K(x_i)))$$

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

## 6.4 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.

## 6.5 Example of regression mixture

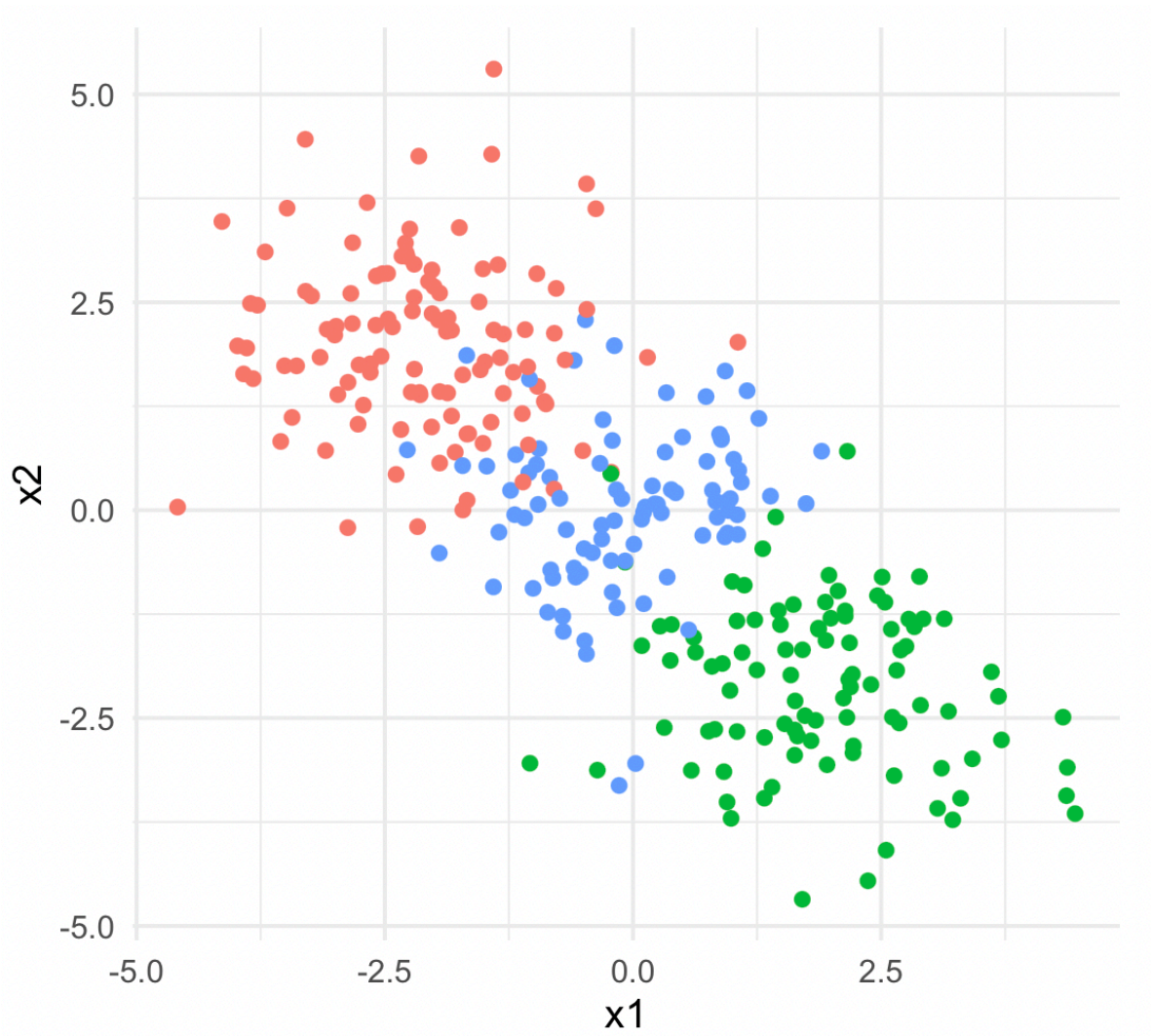
### Training Process

- **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.

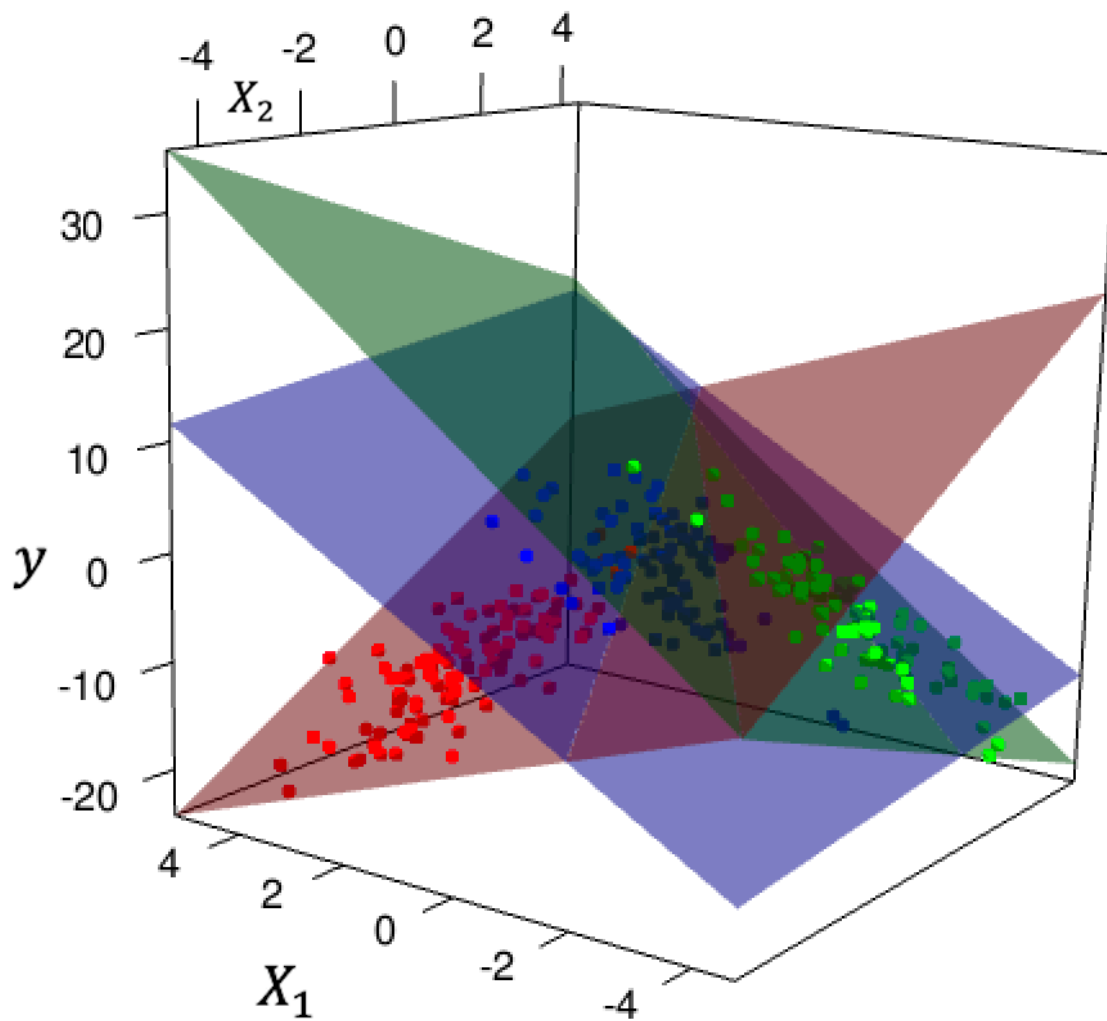
### Inference Phase

- **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.

## 6.6 Example of regression mixture



## Example of regression mixture



## Example of regression mixture

We can model this by allowing the mixing weights and the mixture densities to be input-dependent:

$$p(y_i | x_i, z_i = k, \beta_k, \sigma_k) = \mathcal{N}(y_i | \beta_k^\top x_i, \sigma_k^2 I_2)$$

$$p(z_i | x_i, \theta) = \text{Cat}(z_i | g(\beta^\top x_i))$$

## 7 Singular Value Decomposition

### 7.1 Singular Value Decomposition

#### Eigendecomposition of symmetric matrices

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

$$A = Q\Lambda Q^T$$

#### Singular Value Decomposition

Extend the decomposition to **rectangular matrices**

$$X = USV^T$$

### 7.2 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.
- ...

### 7.3 Existence of the SVD for general matrices

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

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

where  $U^T U = I$  and  $V^T V = I$ .

**In a vector form**

$$X_{n \times d} = \sum_{j=1}^r S_{jj} u_j v_j^T$$

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

### 7.4 Geometrical interpretation

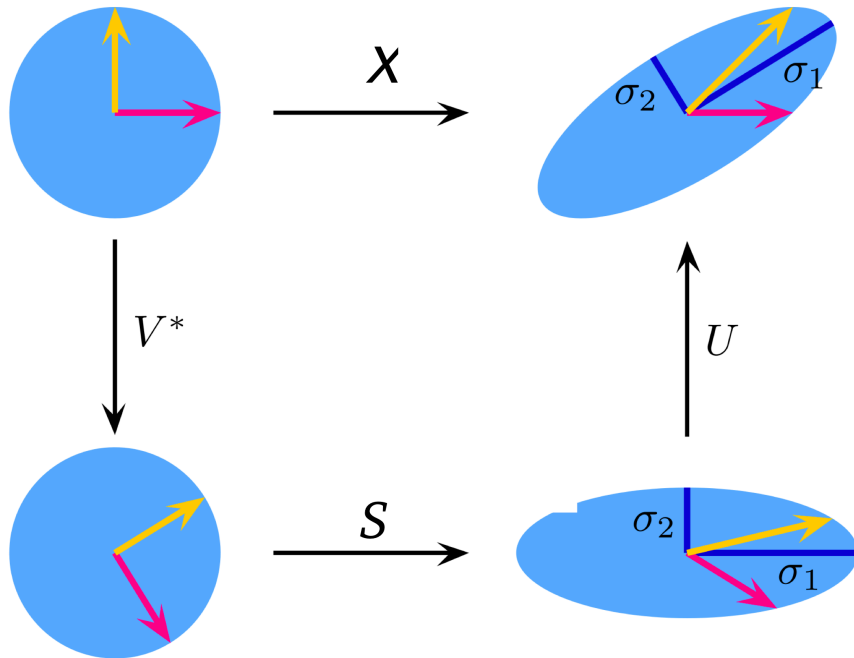
Given any matrix  $X \in \mathbb{R}^{n \times d}$  it defines a linear transformation:

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

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

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

## 7.5 Geometrical interpretation



$$X = U \cdot S \cdot V^*$$

## 7.6 Different versions of SVD

- Full SVD:

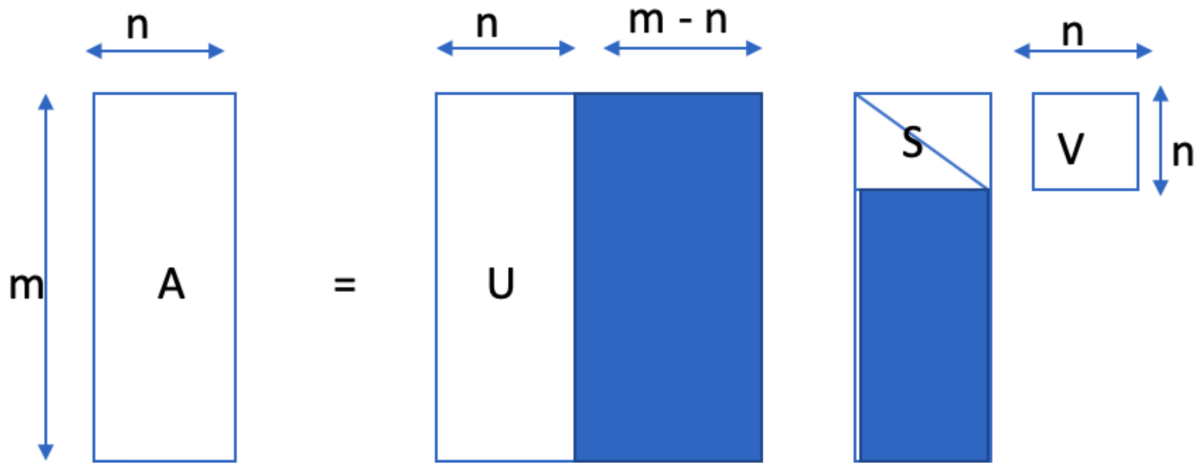
$$X_{n \times d} = U_{n \times n} S_{n \times d} V_{d \times d}^T$$

- Economy sized (thin, compact) SVD:

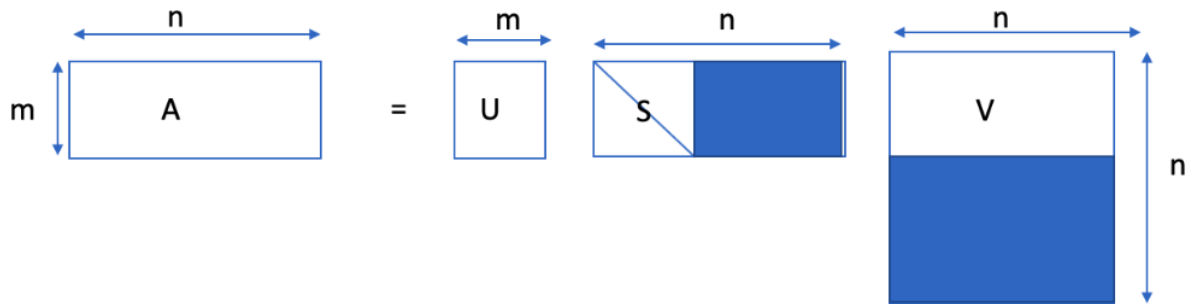
$$X_{n \times d} = U_{n \times r} S_{r \times r} V_{r \times d}^T$$



### 7.7 SVD $n > d$



### 7.8 SVD $n < d$



### 7.9 Existence of the SVD

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

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

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

Define also

$$u_i = \frac{1}{\sigma_i} X v_i \in \mathbb{R}^n,$$

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

## 7.10 Existence of the SVD

### Exercise

It is easy to show that the  $u_1, \dots, u_r$  are orthonormal vectors.

$$u_i^t u_j = \sum_{kl} \frac{v_i^t}{\sigma_i} v_k \Lambda_{kl} v_l^t \frac{v_j}{\sigma_j} = \Lambda_{ij}$$

### Completion if needed

Choose  $u_{r+1}, \dots, u_n \in \mathbb{R}^n$  (through basis completion) such that

$$U = [u_1 \cdots u_n] \in \mathbb{R}^{n \times n}$$

is an orthogonal matrix.

It verifies

$$XV = US,$$

i.e.,

## 7.11 Existence of the SVD

$$X[v_1, \dots, v_r, v_{r+1}, \dots, v_d] = [u_1 \cdots u_r, u_{r+1}, \dots, u_n] \begin{pmatrix} \text{diag}(\sigma_1, \dots, \sigma_r) & 0_{r \times (d-r)} \\ 0_{(n-r) \times r} & 0_{(n-r) \times (d-r)} \end{pmatrix}$$

Two possible cases:

- $1 \leq i \leq r$  :  $Xv_i = \sigma_i u_i$  by construction.
- $i > r$  :  $Xv_i = 0$ , which is due to  $X^T X v_i = C v_i = 0 v_i = 0$ .

Consequently, we have obtained that

$$X = USV^T$$

## 7.12 Properties

The linear application characterized by  $X$  has the following properties:

- $rank(X) = r$  is the number of non zero singular values
- $kernel(X) = span(v_{r+1}, \dots, v_n)$
- $range(X) = span(u_1, \dots, u_r)$

## 7.13 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

## 7.14 A first toy example

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

Original matrix:

```
print(X)
```

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

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

Approximation of rank 2:

```
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
```

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

A basis of the column space:

```
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
```

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

A basis of the kernel:

```
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
```

## 7.15 Illustration of svd in image compression

See the demo of [Tim Baumann](#)

Example borrowed from [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.

```
library(tiff)
library(rasterImage)
```

Loading required package: plotrix

```
peppers<-readTIFF("../Silo-Images/Peppers.tiff")
plot(as.raster(peppers))
```



## 7.16 Illustration of svd in image compression

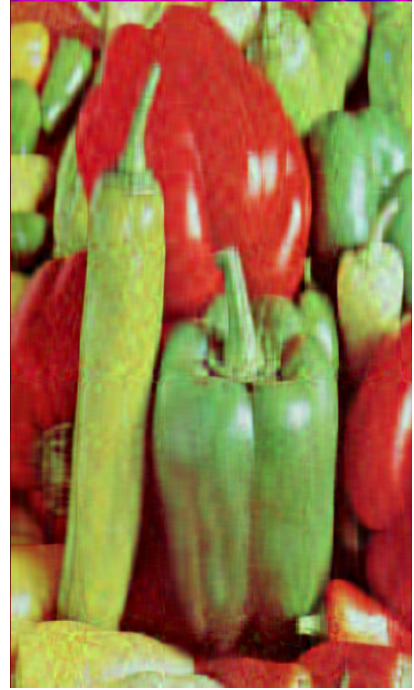
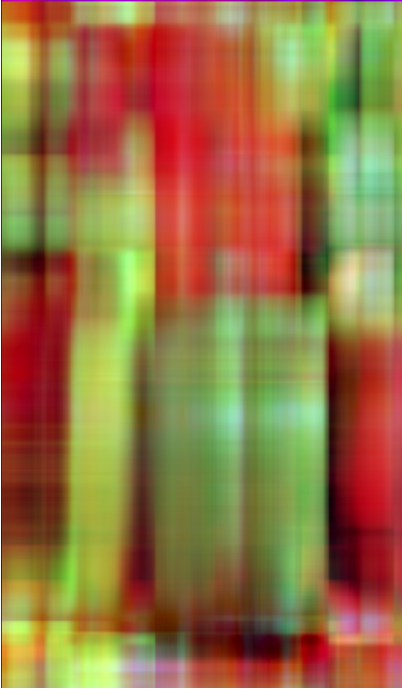
Below the SVD of the three colour intensity matrices, and the view the image that results from using reduced rank versions with rank  $k \in \{5, 30, 100, 300\}$

```
svd_image <- function(im,k){
  s <- svd(im)
  Sigma_k <- diag(s$d[1:k])
  U_k <- s$u[,1:k]
  V_k <- s$v[,1:k]
  im_k <- U_k %*% Sigma_k %*% t(V_k)
  ## the reduced rank SVD produces some intensities <0 and >1.
  # Let's truncate these
  im_k[im_k>1]=1
```

```
im_k[im_k<0]=0
return(im_k)
}

par(mfrow=c(2,2), mar=c(1,1,1,1))

peprssvd<- peppers
for(k in c(4,30,100,300)){
  svds<-list()
  for(ii in 1:3) {
    peprssvd[, ,ii]<-svd_image(peppers[, ,ii],k)
  }
  plot(as.raster(peprssvd))
}
```



## 7.17 Low rank approximation of a matrix $X$

### Frobenius norm

The Frobenius norm of a matrix  $X$  is defined as

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

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

Let

$$\hat{X}_k = \sum_{j=1}^k \sigma_j u_j v_j^T$$

## 7.18 Low rank approximation of a matrix

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

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

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

## 7.19 Proof

We have

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

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

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



## 7.20 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  $i, j \geq 1$

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

- Notice that the  $i$ th biggest singular value of  $X$ ,  $\sigma_i(X)$ , is the first singular value of  $\sigma_1(X - X_{i-1})$
- $\text{rank}(X'_{i-1} + X''_{j-1}) \leq \text{rank}(X_{i+j-2})$  because the columns of each matrix are at most independent and the  $\text{rank}(X'_{i-1} + X''_{j-1}) = i + j - 2$

## 7.21 Proof

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

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

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

## 7.22 Low rank approximation of a matrix and projection

If  $\text{rank}(\hat{X}) = k$ , then we can assume columns  $\hat{X}_i$  of  $\hat{X} \in E_k = \text{span}\{w_1, w_2, \dots, w_k\}$  where  $\{w_1, w_2, \dots, w_k\}$  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  $E_k = \text{span}\{w_1, w_2, \dots, w_k\}$ :

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

where  $\Pi_{E_k} = \sum_{j=1}^k w_j w_j^T$

## 7.23 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  $w_{k+1}, \dots, w_n$  such that  $w_1, \dots, w_n$  form an orthonormal basis of  $\mathbb{R}^n$ . Then,

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

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

## 7.24 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  $E_k$

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

## 7.25 Criterion

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

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

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

## 7.26 Best low rank approximation

$$\hat{X}_k = \Pi_{E_k} X = \sum_{j=1}^k \sigma_j u_j v_j^T = U_{\bullet, 1:k} S_{1:k, 1:k} V_{\bullet, 1:k}^T$$

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

## 7.27 Different views of the approximation

The approximation

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

can be considered in multiple ways:

- approximation of the row
- approximation of the columns

## Notations

If  $X$  is a data table,

- each row  $x_i^T$  is a description of an individual
- each column  $X_j$  is variable describing  $n$  individuals

## 7.28 Rows approximation (projection of the individuals)

Transposing the matrix the best low rank approximation becomes

$$\hat{X}_k^T = \Pi_{F_k} X^T = \sum_{j=1}^k \sigma_j v_j u_j^T = V_{\bullet,1:k} S_{1:k,1:k} U_{\bullet,1:k}^T$$

where  $F_k = \text{span}\{v_1, \dots, v_k\}$

## 7.29 The approximation error

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

Each row  $x_i$  is approximated by

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

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

## 7.30 Projection of the variables

$\Pi_{E_k}$  is the projection matrix on  $E = \text{span}(u_1, \dots, u_k)$

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

and

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

## 7.31 k first principal components

$$C_{1:n,1:k} = U_{E_k} S_{1:k,1:k}$$

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

The principal component are the coordinates of the projection of the rows of  $X$  on  $F_k$ :

$$C_{1:n,1:k} = X V_{1:d,1:k}$$

### 7.32 Percentage of information

We have  $C_{\bullet,1:k}^T C_{\bullet,1:k} = S_{1:k,1:k}^2 = \text{diag}(\sigma_1^2, \dots, \sigma_k^2)$ , thus

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

and

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

### 7.33 Correlations

$$\widehat{\text{cor}}(X_{\bullet,j}, C_{\bullet,k}) = \frac{X_{\bullet,j}^T C_{\bullet,k}}{\|X_{\bullet,j}\| \|C_{\bullet,k}\|} = \cos(X_{\bullet,j}, \widehat{C_{\bullet,k}})$$

### 7.34 Duality

It is easy to show that

- the columns of  $V$  are the eigenvector of  $X^T X$
- the columns of  $U$  are the eigenvector of  $X X^T$

Thus the principal component of  $X^T X$  are the eigenvectors of  $X X^T$  and vice-versa

## 8 Multi Dimensional Scaling

### 8.1 Dimensionality reduction and manifold learning

#### Goal

Given pairwise

- dissimilarities  $\delta_{ij}$  or
- high-dimensional input data  $\{x_i\}_{i=1,\dots,n}$

reconstruct a lower dimensional map with embedded data  $\{y_i\}_{i=1,\dots,n}$

## Principle

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

## 8.2 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.

## 8.3 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:

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

## 8.4 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.

## 8.5 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

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

## Distance and data table (continued)}

- Let  $J$  be the centering matrix,

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

- Centering of  $X$  in columns

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

- Centering of  $X$  in rows

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

## 8.6 Distance and data table (continued)}

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

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

## 8.7 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.

## 8.8 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**

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



## 8.9 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!

## 8.10 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  $d_{ij}$

## 8.11 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

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

## 8.12 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.

### 8.13 Normalized error function

#### Normalized Stress

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

#### A general form

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

### 8.14 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

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

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

- $f$  is a monotone function that preserves order relations
- if  $\delta_{ij} < \delta_{kl}$  then  $f(\delta_{ij}) < f(\delta_{kl})$

## 8.15 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

$$d_{ij} = d(\mathbf{y}_i, \mathbf{y}_j) \tag{1}$$

$$= \|\mathbf{y}_i - \mathbf{y}_j\|_2 \tag{2}$$

## 8.16 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

## 8.17 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

## 8.18 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)
```

## 8.19 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.  $\frac{\partial \|y_i - y_j\|}{\partial y_{ik}} = \frac{\partial d_{ij}}{\partial d_{ij}^2} \frac{\partial d_{ij}^2}{\partial y_{ik}} = \frac{y_{ik} - y_{jk}}{d_{ij}}$

## 8.20 Sammon gradient descent

At iteration  $q$

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

## 8.21 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.

```
Loading required package: colorspace
```

```
Loading required package: e1071
```

```
Attaching package: 'smacof'
```

```
The following object is masked from 'package:base':
```

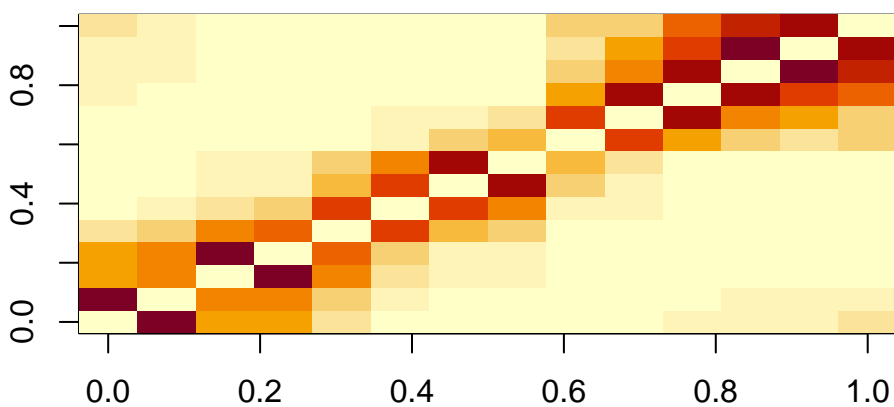
```
transform
```

```
Loading required package: photobiology
```

```
Documentation at https://docs.r4photobiology.info/
```

## 8.22 Example Ekman colors

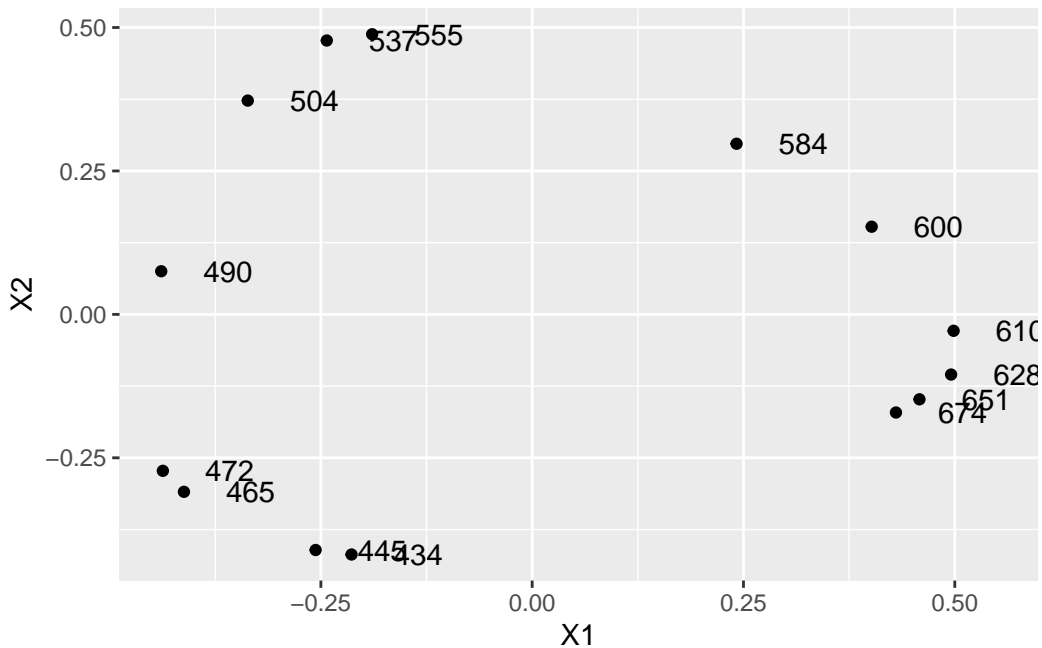
```
image(as.matrix(ekman))
```



## 8.23 Example Classical scaling

```
ekmanD <- sim2diss(ekman, method = "confusion", to.dist = TRUE)
Y<-data.frame(cmdscale(ekmanD))
ggplot(data=Y,aes(x=X1,y=X2,label=rownames(Y)),color=as.double(rownames(Y)))+geom_point()+geom
```

Warning: Computation failed in `stat\_color()`.  
Caused by error in `color\_of.numeric()`:  
! all(ifelse(is.na(x), Inf, x) > 0) is not TRUE



## 9 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

## 9.1 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  $x_i$ :

## 9.2 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

### 9.3 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).

### 9.4 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.

### 9.5 Gradient descent

#### Gradient

$$\frac{\partial C}{\partial y_i} = 2 \sum_j (p_{j|i} - q_{j|i} + p_{i|j} - q_{i|j})(y_i - y_j).$$



## Interpretation

The gradient may be interpreted as the resultant force created by a *set of springs* between the map point  $y_i$  and all other map points  $y_j$ .

- force exerted by the spring between  $y_i$  and  $y_j$  is proportional to its length
- also proportional to its stiffness  $(p_{j|i} - q_{j|i} + p_{i|j} + q_{i|j})$

## Initialisation

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

## 9.6 Gradient in practice

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

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

where  $y^{(t)}$  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

## 9.7 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*.

## 9.8 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**  $q_{ij}$

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

## 9.9 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.

- $p_{ij} = \frac{p_{ji} + p_{ij}}{2n}$ . 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

## 9.10 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

## 9.11 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

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}}$$

**Gradient**

$$\frac{\partial C}{\partial y_i} = 4 \sum_j (p_{ij} - q_{ij})(y_i - y_j)(1 + \|y_i - y_j\|^2)^{-1}. \quad (3)$$

## 9.12 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

### 9.13 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

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

### 9.14 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

### 9.15 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.

## 10 Word Embedding

## 11 Word embedding

### 11.1 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.

### 11.2 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)*

### 11.3 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

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

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

## 11.4 Cosine for measuring similarity

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

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

## 11.5 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.

## 11.6 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  $idf_t$ :

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

	As You Like It	Twelfth Night	Julius Caesar	Henry V
<b>battle</b>	0.074	0	0.22	0.28
<b>good</b>	0	0	0	0
<b>fool</b>	0.019	0.021	0.0036	0.0083
<b>wit</b>	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.*

## 11.7 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:

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

## 11.8 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)

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

	<b>computer</b>	<b>data</b>	<b>result</b>	<b>pie</b>	<b>sugar</b>
<b>cherry</b>	0	0	0	4.38	3.30
<b>strawberry</b>	0	0	0	4.10	5.51
<b>digital</b>	0.18	0.01	0	0	0
<b>information</b>	0.02	0.09	0.28	0	0

*PPMI matrix showing the association between words and context words*

## 11.9 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.

## 11.10 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.

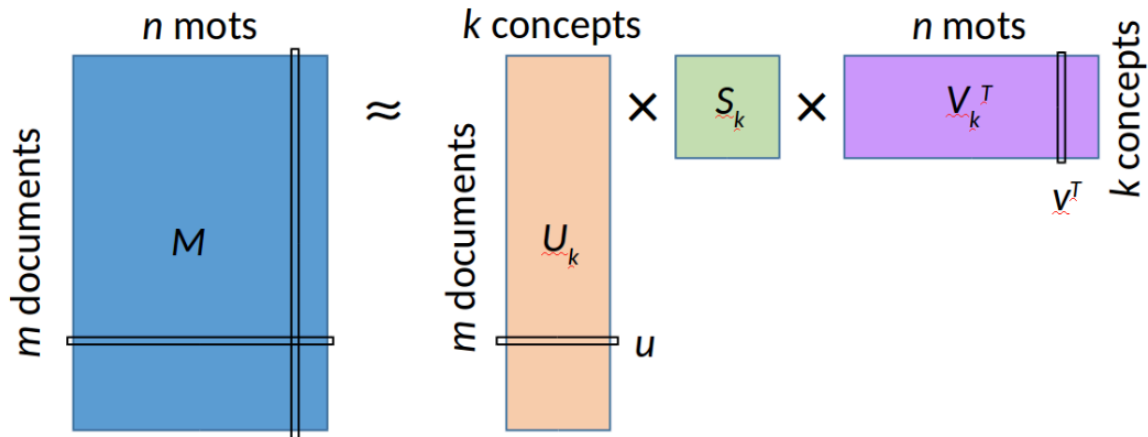
### Principle

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



## 11.11 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  $U_k S_k V_k^T$



## 11.12 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

$$U_k S_k \mathbf{v}^T$$

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

$$\mathbf{u} S_k V_k^T$$

## 11.13 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.

## 11.14 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  $w_{i,j} \sim \text{Cat}(\varphi_{z_{i,j}})$ .

### 11.15 Latent Dirichlet Allocation (LDA) in summary

$$\varphi_{k=1\dots K} \sim \text{Dirichlet}_V(\beta) \quad (4)$$

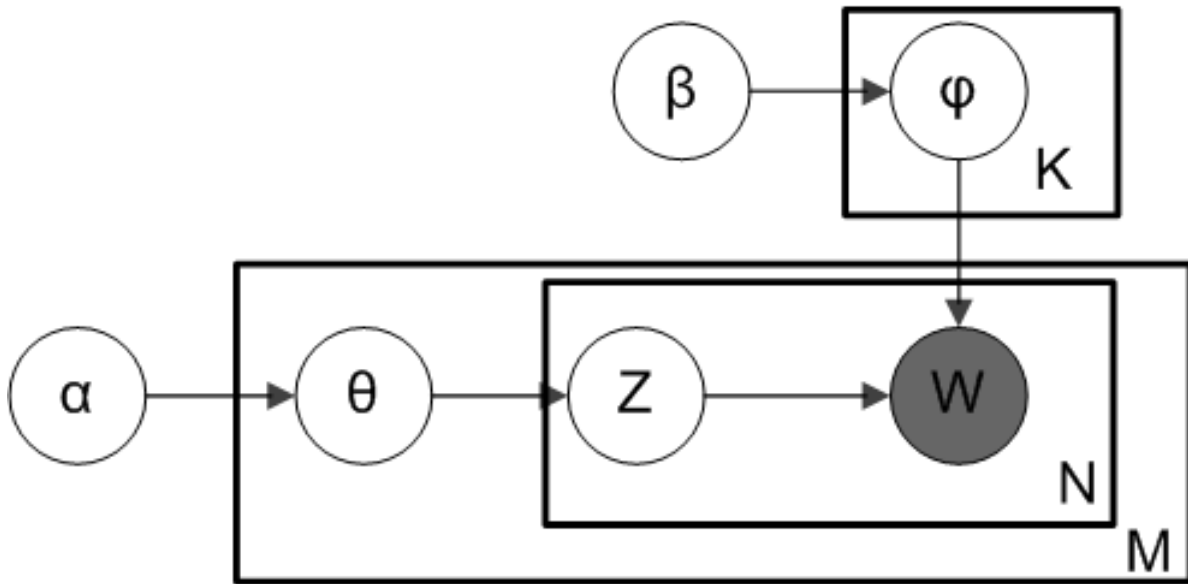
$$\theta_{d=1\dots M} \sim \text{Dirichlet}_K(\alpha) \quad (5)$$

$$z_{d=1\dots M, w=1\dots N_d} \sim \text{Categorical}_K(\theta_d) \quad (6)$$

$$w_{d=1\dots M, w=1\dots N_d} \sim \text{Categorical}_V(\varphi_{z_{dw}}) \quad (7)$$

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

### 11.16 Latent Dirichlet Allocation (LDA) estimation



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

### 11.17 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

## 11.18 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

- Topic Modeling: Organizing large text corpora (news articles, research papers).
- Recommendation Systems: Suggesting articles based on topic similarity.
- Text Classification: Assigning labels based on topic distributions.
- Social Media Analysis: Identifying trending topics in posts.

## 11.19 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)

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

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

### 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:

$$\text{Similarity}(w, c) = c^T w$$

## 11.21 From dot product to logistic regression

### Logit

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

## 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  $k$ ).

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

### 11.22 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 ...

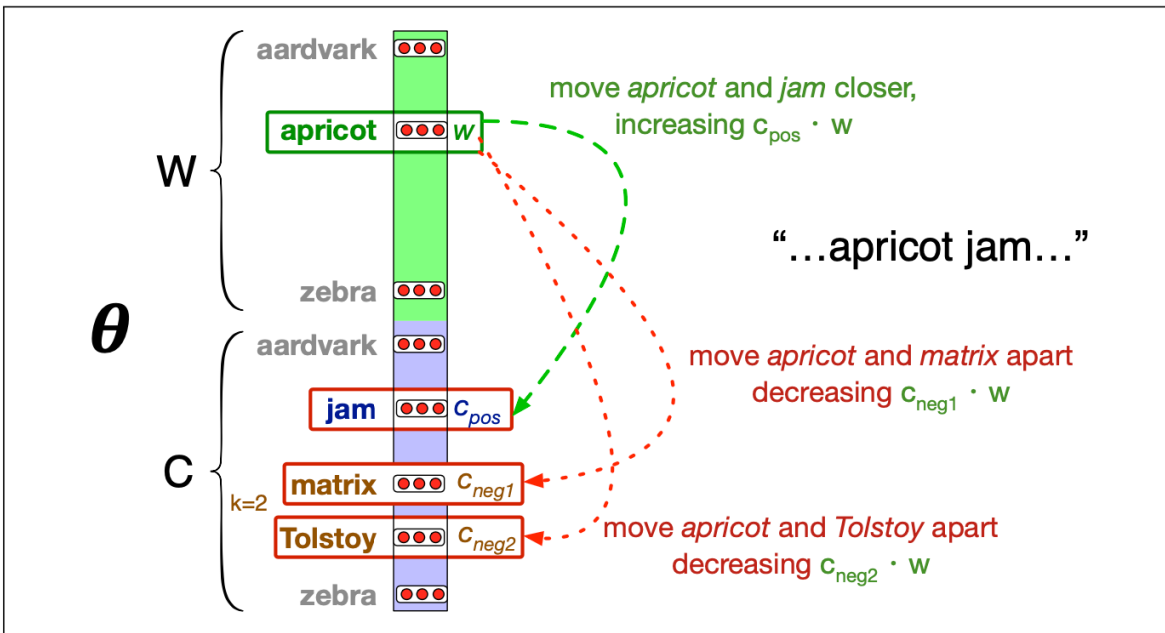
#### positive examples +

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

#### negative examples -

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

### 11.23 Skip-gram model learns two separate embeddings



### 11.24 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.

## 12 Transformers

### 12.1 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**

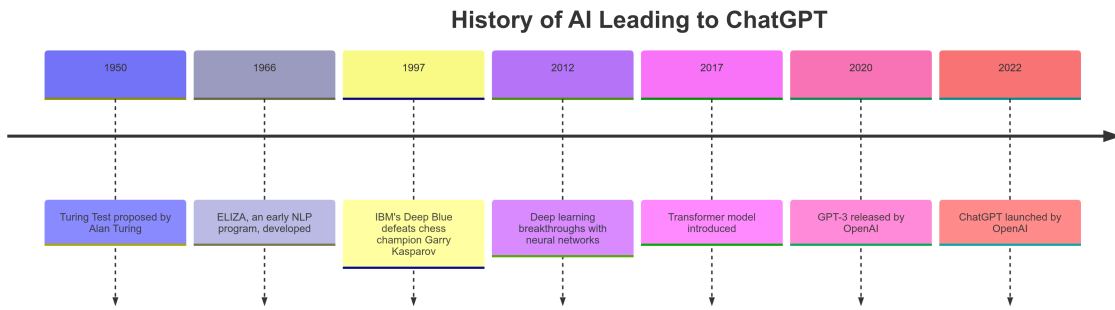


Figure 8: History of AI leading to chatGPT

## 12.2 The nature of transformers

A classical transformer is a function

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

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

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

Figure 9: 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

- $w_{1:n}$  is usually coded a matrix of one-hot vectors with  $w_{1:n} \in \mathbb{R}^{n \times |V|}$ .
- For classical chatbot  $p_j = P(w_{n+1} = w_{(j)} | w_{1:n})$  where  $w_{(j)}$  is the  $j$ th element of  $\mathcal{V}$ .

## 12.3 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

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

## 12.4 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  $(s_j)$ :

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

## 12.5 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

## 12.6 Initial embedding

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

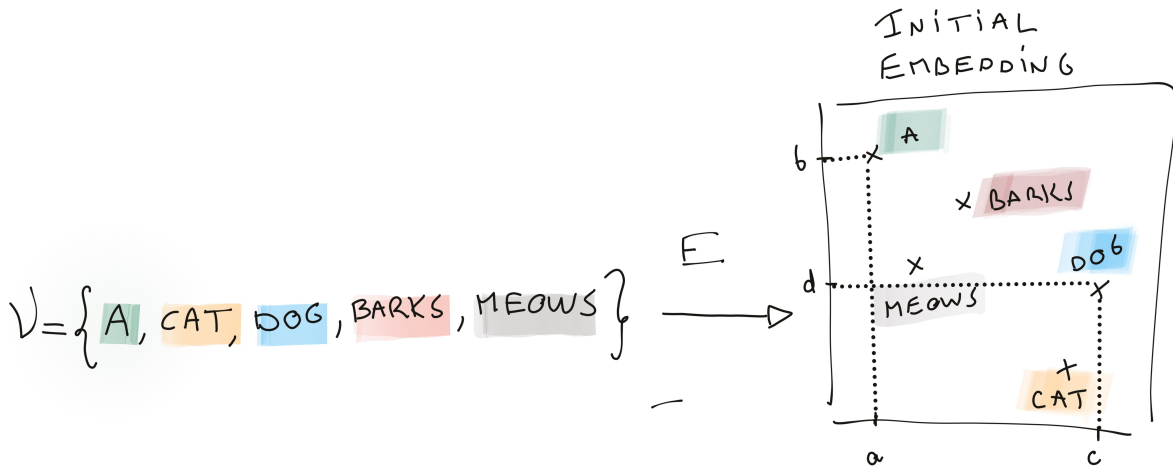


Figure 10: Initial embedding

## 12.7 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  $d \times d$ .

## 12.8 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  $V$

## 12.9 Illustration of self-Attention

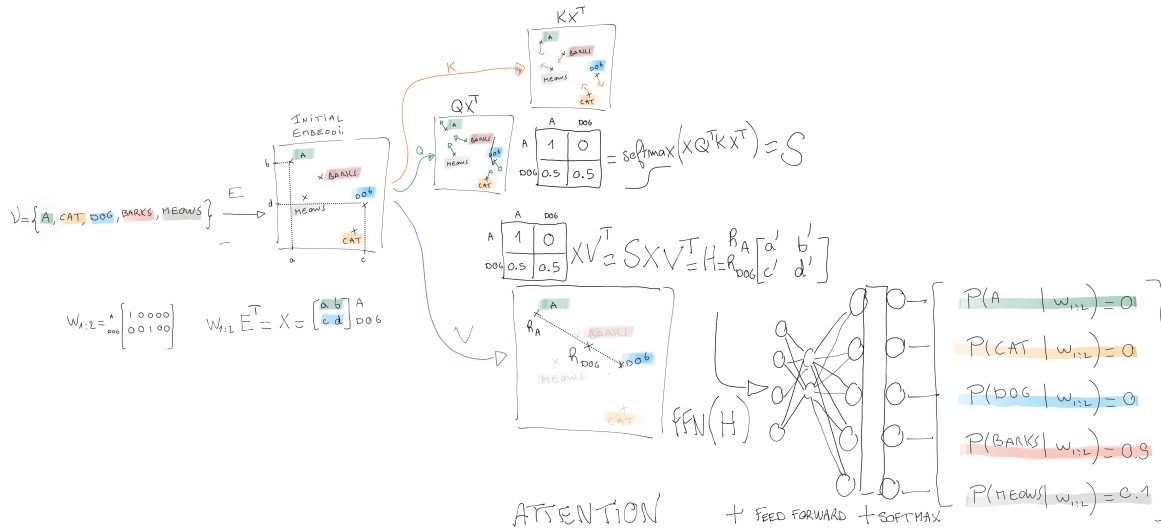


Figure 11: Self-attention block in auto-encoder

## 12.10 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

## 12.11 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

## 12.12 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

## 12.13 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.

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

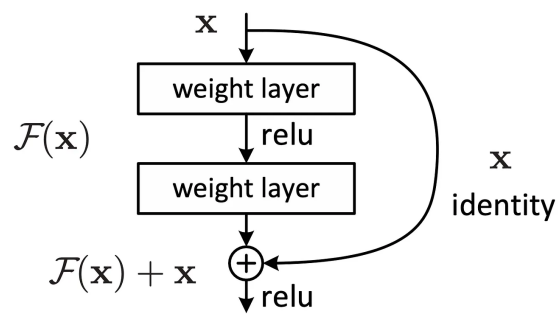


Figure 12: Skip connection

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

## 12.15 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.

## 12.16 Fine tuning

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

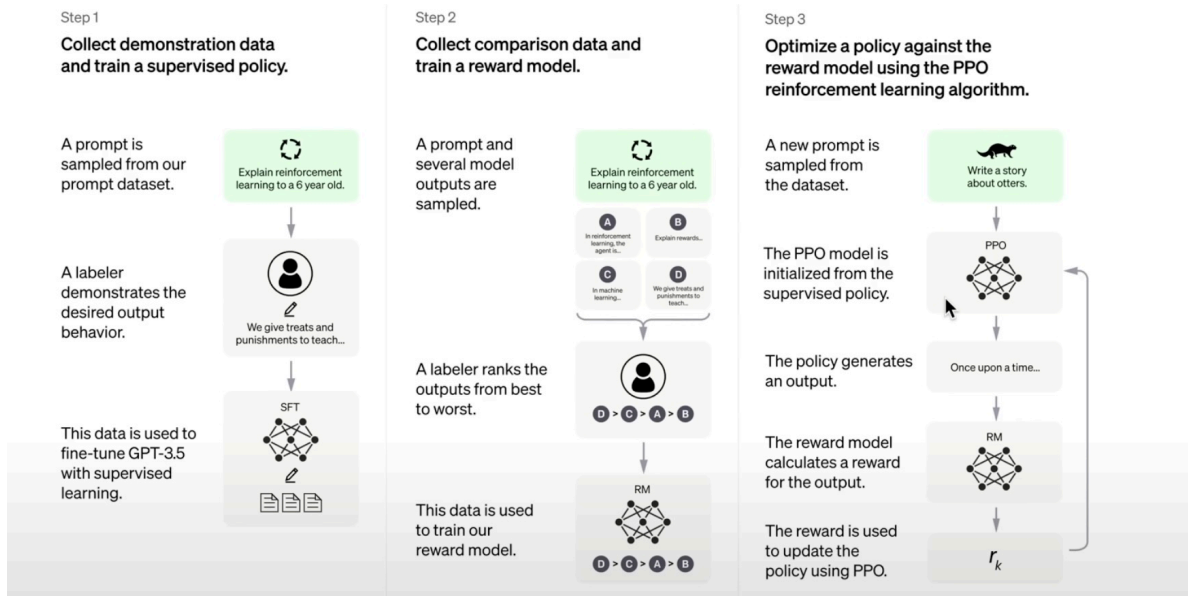


Figure 13: ChatGPT fine tuning (from A. Karpathy lecture)

## 13 Exercises

### 13.1 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  $AX = B$ , where  $A$  is invertible and partitioned as described above.
2. Write the equation using the partitioned form of matrix  $A$ :

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix}$$

3. Apply block matrix operations to rewrite the equation:

$$A_{11}X_1 + A_{12}X_2 = B_1 \quad (8)$$

$$A_{21}X_1 + A_{22}X_2 = B_2 \quad (9)$$

4. Solve the first equation for  $X_1$ :

$$X_1 = A_{11}^{-1}(B_1 - A_{12}X_2)$$

5. Substitute this value of  $X_1$  into the second equation:

$$A_{21}A_{11}^{-1}(B_1 - A_{12}X_2) + A_{22}X_2 = B_2$$

6. Simplify the equation:

$$SX_2 = B_2 - A_{21}A_{11}^{-1}B_1$$

7. We can see that the coefficient matrix for  $X_2$  is the Schur complement,  $S = A_{22} - A_{21}A_{11}^{-1}A_{12}$

8. Therefore,  $X_2$  can be calculated as:

$$X_2 = S^{-1}(B_2 - A_{21}A_{11}^{-1}B_1)$$

9. Finally, we can substitute the values of  $X_1$  and  $X_2$  to obtain the solution vector  $X$ :

$$\begin{aligned} X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} &= \begin{bmatrix} A_{11}^{-1}(B_1 - A_{12}S^{-1}(B_2 - A_{21}A_{11}^{-1}B_1)) \\ S^{-1}(B_2 - A_{21}A_{11}^{-1}B_1) \end{bmatrix} = \\ &= \begin{bmatrix} A_{11}^{-1} + A_{11}^{-1}A_{12}S^{-1}A_{21}A_{11}^{-1} & -A_{11}^{-1}A_{12}S^{-1} \\ -S^{-1}A_{21}A_{11}^{-1} & S^{-1} \end{bmatrix} \begin{bmatrix} B_1 \\ B_2 \end{bmatrix} \end{aligned}$$

This demonstrates the Schur complement lemma and its application in solving the matrix equation  $AX = B$ .

Memotechnic : Take the other block in the diagonal and then turn clockwise to find the product of the three other block: for  $A_{11}$  we get  $A_{11} - A_{12}A_{22}^{-1}A_{21}$

## 13.2 Conditional Gaussian

Demonstrate the form of the distribution of  $x_1|x_2$  when both vectors are gaussian.

We have the joint gaussian distribution of  $\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ :

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \sim \mathcal{N}(\mu, \Sigma)$$

We can proceed in three step:

1. Compute  $\Sigma^{-1} = \Lambda$  the concentration matrix
2. Write the joint distribution into a sum of marginal in  $x_2$  and conditonal in  $x_1$
3. Derive the conditional distribution.

### Concentration matrix

The joint covariance matrix is partitioned as:

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

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

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

where the blocks of the concentration matrix are given by:

$$\Lambda_{11} = (\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21})^{-1}, \quad \Lambda_{12} = -\Lambda_{11}\Sigma_{12}\Sigma_{22}^{-1},$$

$$\Lambda_{21} = -\Sigma_{22}^{-1}\Sigma_{21}\Lambda_{11}, \quad \Lambda_{22} = \Sigma_{22}^{-1} + \Sigma_{22}^{-1}\Sigma_{21}\Lambda_{11}\Sigma_{12}\Sigma_{22}^{-1}.$$

Here,  $\Lambda_{11}$  is obtained using the Schur complement of  $\Sigma_{11}$  in

$\Sigma$ :

$$\Lambda_{11} = (\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21})^{-1}.$$

## Decompose the joint distribution

The quadratic form of the joint distribution decomposes as follow

$$Q = (x_1 - \mu_1)^T \Lambda_{11} (x_1 - \mu_1) + 2(x_1 - \mu_1)^T \Lambda_{12} (x_2 - \mu_2) + (x_2 - \mu_2)^T \Lambda_{22} (x_2 - \mu_2)$$

Replacing the  $\Lambda_{ij}$  with their expression in function of the  $\Sigma_{ij}$ , we eventually get

$$Q = (x_1 - (\mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2)))^T \Sigma_{1|2}^{-1} (x_1 - (\mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2))) + \quad (10)$$

$$(x_2 - \mu_2)^T \Sigma_{22}^{-1} (x_2 - \mu_2) \quad (11)$$

where

$$\Sigma_{1|2} = (\Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21})^{-1} = \Lambda_{11}$$

## Derive the conditional distribution

The joint distribution of  $(x_1, x_2) \sim N(\mu, \Sigma)$ , where

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix},$$

can be expressed as a product of the marginal distribution of  $x_2$  and the conditional distribution of  $x_1$  given  $x_2$ .

The joint log-probability is:

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

Expanding each term: 1. The marginal distribution of  $x_2$ :

$$p(x_2) \sim N(\mu_2, \Sigma_{22}),$$

with the log-density:

$$\log p(x_2) = -\frac{1}{2} (x_2 - \mu_2)^T \Sigma_{22}^{-1} (x_2 - \mu_2) - \frac{1}{2} \log |\Sigma_{22}| - \frac{d_2}{2} \log(2\pi).$$

2. The conditional distribution of  $x_1$  given  $x_2$ :

$$p(x_1|x_2) \sim N(\mu_{1|2}, \Sigma_{1|2}),$$

where

$$\mu_{1|2} = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(x_2 - \mu_2), \quad \Sigma_{1|2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}.$$

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:

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

### 13.3 Stochastic Gradient and linear regression

The Mean Squared Error can be expressed as a sum of the available sample  $(y_i, x_i)_i$ :

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

#### The gradient

is expressed as

$$\nabla Q_i(w) = -2x_i(y_i - w^t x_i)$$

## Online regression code

```
library(ggplot2)
library(patchwork) # Pour afficher les graphes côte à côte
```

Attaching package: 'patchwork'

The following object is masked from 'package:MASS':

area

The following object is masked from 'package:cowplot':

align\_plots

```
# Fonction de régression linéaire séquentielle
regression.online <- function(X, Y, max.iteration = 300) {
  p <- ncol(X)
  n <- nrow(X)
  X <- cbind(1, X) # Ajout de l'intercept

  shuffling <- sample(1:n, n)
  X <- X[shuffling, ]
  Y <- Y[shuffling]
  Q<-rep(0,max.iteration)
  W <- rep(0, p + 1) # Initialisation des coefficients
  for (i in 1:max.iteration) {
    idx <- (i - 1) %% n + 1
    x <- drop(X[idx,]) # Sélection d'un point (matrice 1x(p+1))
    y <- Y[idx]
    W <- W + (1 / i) * x * (y - drop(rbind(W) %*% cbind(x))) # Mise à jour du vecteur de po
    Q[i]<-mean((Y - X %*% cbind(W))^2)
  }

  return(list(W=W,Q=Q))
}
```

## Comparison with batch approach

```
# Génération des données
set.seed(123)
x <- rnorm(200, sd = 2)
y <- 1 + 3 * x + rnorm(200, sd = 1)

# Ajustement avec lm
lm_model <- lm(y ~ x)
lm_coef <- coef(lm_model) # Coefficients de lm

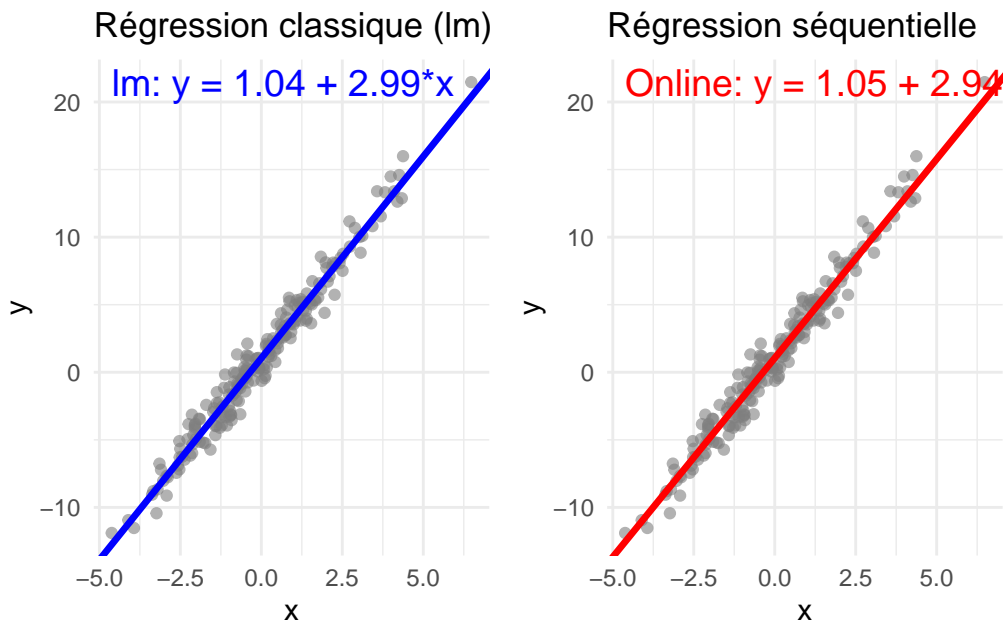
# Ajustement avec la régression séquentielle
online_res <- regression.online(matrix(x, ncol = 1), y)
online_coef <- online_res$W

# Création d'un dataframe pour ggplot
df <- data.frame(x = x, y = y)

# Graphique 1 : Régression classique avec lm
p1 <- ggplot(df, aes(x = x, y = y)) +
  geom_point(color = "gray50", alpha = 0.6) +
  geom_abline(intercept = lm_coef[1], slope = lm_coef[2], color = "blue", lwd = 1.2) +
  annotate("text", x = min(x), y = max(y),
          label = sprintf("lm: y = %.2f + %.2f*x", lm_coef[1], lm_coef[2]),
          hjust = 0, color = "blue", size = 5) +
  labs(title = "Régression classique (lm)",
       x = "x", y = "y") +
  theme_minimal()

# Graphique 2 : Régression en ligne
p2 <- ggplot(df, aes(x = x, y = y)) +
  geom_point(color = "gray50", alpha = 0.6) +
  geom_abline(intercept = online_coef[1], slope = online_coef[2], color = "red", lwd = 1.2) +
  annotate("text", x = min(x), y = max(y),
          label = sprintf("Online: y = %.2f + %.2f*x", online_coef[1], online_coef[2]),
          hjust = 0, color = "red", size = 5) +
  labs(title = "Régression séquentielle",
       x = "x", y = "y") +
  theme_minimal()

# Affichage côte à côte
p1 + p2
```



### MSE minimisation

```
library(ggplot2)

# Données pour la courbe de la régression online
df_online <- data.frame(
  MSE = online_res$Q,
  Iteration = 1:length(online_res$Q)
)

# MSE de la régression classique
mse_lm <- mean(lm_model$residuals^2)

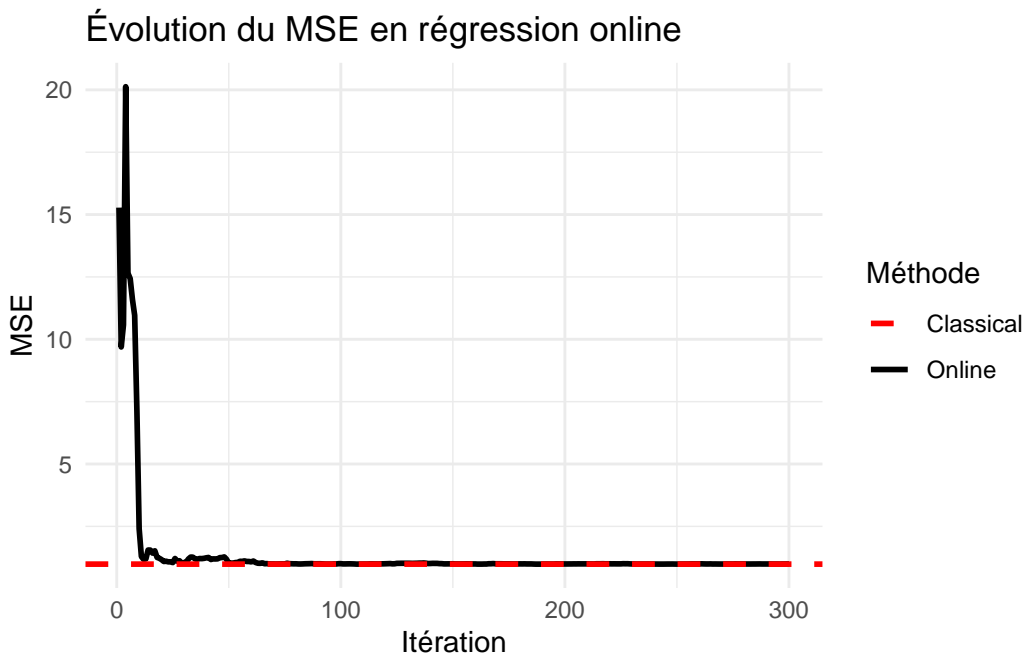
# Graphique avec ggplot2
ggplot(df_online, aes(x = Iteration, y = MSE)) +
  geom_line(color = "black", size = 1, aes(linetype = "Online")) + # Courbe Online en noir
  geom_hline(aes(yintercept = mse_lm, linetype = "Classical"), color = "red", size = 1) + #
  scale_linetype_manual(name = "Méthode", values = c("Online" = "solid", "Classical" = "dash")) +
  labs(
    title = "Évolution du MSE en régression online",
    x = "Itération",
```

```

y = "MSE"
) +
theme_minimal()

```

Warning: Using `size` aesthetic for lines was deprecated in ggplot2 3.4.0.  
 i Please use `linewidth` instead.



### 13.4 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:

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

Now, let our two Normal distributions be  $\mathcal{N}(\mu_p, \Sigma_p)$  and  $\mathcal{N}(\mu_q, \Sigma_q)$ , both  $k$  dimensional.



$$D_{KL}(p||q) = \mathbb{E}_p [\log(p) - \log(q)] = \mathbb{E}_p \left[ \frac{1}{2} \log \frac{|\Sigma_q|}{|\Sigma_p|} - \frac{1}{2}(\mathbf{x} - \mu_p)^T \Sigma_p^{-1} (\mathbf{x} - \mu_p) + \frac{1}{2}(\mathbf{x} - \mu_q)^T \Sigma_q^{-1} (\mathbf{x} - \mu_q) \right] = \frac{1}{2} \log \frac{|\Sigma_q|}{|\Sigma_p|} + \frac{1}{2} \mathbb{E}_p [(\mathbf{x} - \mu_q)^T \Sigma_q^{-1} (\mathbf{x} - \mu_q)] - \frac{1}{2} \mathbb{E}_p [(\mathbf{x} - \mu_p)^T \Sigma_p^{-1} (\mathbf{x} - \mu_p)] \quad (12)$$

Now, since  $(\mathbf{x} - \mu_p)^T \Sigma_p^{-1} (\mathbf{x} - \mu_p)$  in the second term  $\in \mathbb{R}$ , we can write it as  $tr \{ (\mathbf{x} - \mu_p)^T \Sigma_p^{-1} (\mathbf{x} - \mu_p) \}$ , 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  $tr \{ (\mathbf{x} - \mu_p)(\mathbf{x} - \mu_p)^T \Sigma_p^{-1} \}$ .

The second term now is,

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

The expectation and trace can be interchanged to get,

$$= \frac{1}{2} tr \{ \mathbb{E}_p [(\mathbf{x} - \mu_p)(\mathbf{x} - \mu_p)^T \Sigma_p^{-1}] \} = \frac{1}{2} tr \{ \mathbb{E}_p [(\mathbf{x} - \mu_p)(\mathbf{x} - \mu_p)^T] \Sigma_p^{-1} \} \quad (13)$$

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

$$= \frac{1}{2} tr \{ \Sigma_p \Sigma_p^{-1} \} = \frac{1}{2} tr \{ I_k \} = \frac{k}{2} \quad (14)$$

We can simplify the third term

$$\begin{aligned} \mathbb{E}_p [(\mathbf{x} - \mu_q)^T \Sigma_q^{-1} (\mathbf{x} - \mu_q)] &= \mathbb{E}_p [(\mathbf{x} - \mu_p + \mu_p - \mu_q)^T \Sigma_q^{-1} (\mathbf{x} - \mu_p + \mu_p - \mu_q)] \\ &= \mathbb{E}_p [(\mu_p - \mu_q)^T \Sigma_q^{-1} (\mu_p - \mu_q) + tr \{ \Sigma_q^{-1} \Sigma_p \} + 2(\mu_p - \mu_q)^T \Sigma_q^{-1} (\mu_p - \mu_q)] \\ &= (\mu_p - \mu_q)^T \Sigma_q^{-1} (\mu_p - \mu_q) + tr \{ \Sigma_q^{-1} \Sigma_p \} \end{aligned} \quad (15)$$

Combining all this we get,

## 13.5 Kullback–Leibler divergence between two multivariate Gaussian distributions

**Compact form**

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

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

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