Unweaving massive data with sparse coding

Djork-Arné Clevert

Institute of Bioinformatics, Johannes Kepler University Linz

Boston, July 11th, 2014

Objectives

- What is a sparse coding?
- What are some common sparse coding models and what are their fields of application?
- What is the role of assumptions in the data generation process?
- How can domain knowledge be incorporated in order to build an appropriate model?

Sparse coding



Pergamon

PH: \$8042-6989(97)00169-7

Pinian Brz., Vol. 37, No. 23, pp. J311-3325, 1997 S 1991 Elseviar brianto I.id. All rightn material Printed in Secal Britain 9042-6488/97 S 17.341 + 3104

Sparse Coding with an Overcomplete Basis Set: A Strategy Employed by V1?

BRUNG A. OLSHAUSEN.; DAVID J. FIELD; Reveived 16 July 1996; in revised form 24 December 1996.

Given a large set of input patterns, sparse coding attempt to automatically find a small number of representative patterns which reproduce the original input patterns. Motivation

Sparse coding and linear generative models

- Some sources 'drive' the model and produce an output distribution which should best match the observed data distribution
- Observations x are random variables whose distribution depends on model parameters $\boldsymbol{\theta}$
- Nice advantages of generative models are:
 - they select models using sound model selection techniques such as maximum likelihood or maximum a posterior,
 - signal-to-noise ratios can be computed,
 - they can be compared with each other via the likelihood or posterior,
 - they produce a global model to explain all data.

Some applications

Factor analysis for

- quantitative gene expression analysis
 - Factor Analysis for Robust Microarray Summarization (FARMS)
- discovering copy number variations in array data
 - copy number FARMS (cn.FARMS)
- identifying bicluster in a data matrix
 - Factor Analysis for Bicluster Acquisition (FABIA)

Microarray





Facts and assumptions

- Gene measured by different probes
- Goal: summarize probe intensities to an expression value
- Noise-free probes are correlated
- Replicate probe intensities are Gaussian distributed

20	1	1						15		
Probe 1								مرجع	 	
	Probe 2	10			1	1	100	1	1.1	
	ġ.	Probe 3	and a	- 1	a ser			and a second		
		a.	Probe 4		100			10		8°**
X	¥19.	A.	8. ⁻	Probe 5	1		1	~	~~~	
196	<u>.</u>		ŵ.		Probe 6				- ^	6°**
						Probe 7	<i></i>	1	<u> </u>	
					5	17	Probe 8	1	<u></u>	
					247 22	- 2 55		Probe 9		
	ingent Anter		ave.	inge. Tado I	1994). 1986)	1738) 1. at 4	and a Cantar		Probe 10	
البتبا	- 10 - 645 -	Ļ				لتبتبا				[¹⁰⁰⁰]

ロト (個) (主) (主) した のへの

The probe level data



FARMS: The idea



ロト (個) (注) (主) (日) のへの

Factor analysis

$$x = \lambda z + \epsilon = \sum_{i=1}^{p} \lambda_i z_i^T + \epsilon$$

- linear generative model
- $x \in \mathbb{R}^n$ are the observations
- $z \in \mathbb{R}^p$ are independent $\mathcal{N}(0,1)$ Gaussian factor
- $\boldsymbol{\epsilon}$ is independent $\mathscr{N}(\mathbf{0},\Psi)$ Gaussian noise
- x is Gaussian with $P(x) = \int P(z) P(x | z) dx = \mathcal{N}(0, \lambda \lambda^{t} + \Psi)$
- where $\mathbf{\lambda} \in \mathbb{R}^{n imes p}$ is the factor loading matrix,
- and the noise covariance matrix $\Psi \in \mathbb{R}^{p \times p}$ is diagonal

MA-plot MAS5.0 vs. FA



The MA plot shows log fold change as a function of mean log expression level.

Facts and assumptions II

- Gene measured by different probes
- Goal: summarize probe intensities to an expression value
- Noise-free probes are positively correlated
 - Variable probe qualities
 - High quality probes are linear dependent
- Replicate probe intensities are Gaussian distributed



Higher mRNA concentration \rightarrow larger intensities

FARMS: Bayes framework

FARMS

Posterior

$$p(\mathbf{\lambda}, \mathbf{\Psi} | \{x\}) \propto p(\{x\} | \mathbf{\lambda}, \mathbf{\Psi}) p(\mathbf{\lambda})$$

Prior knowledge

- Positive λ ensure positive probe correlation
- Most genes show no or small signal (large signals are of interest in a study)

Rectified Gaussian



$$\begin{split} \lambda_{j} &= \max\{y_{j}, 0\} \text{ with } \\ y_{j} &\sim \mathcal{N}\left(\mu_{\lambda}, \sigma_{\lambda}\right) \end{split}$$

FARMS: EM updates

FARMS

E-step:

$$\mathsf{E}_{\mathbf{z}_i \mid \mathbf{x}_i}\left(z_i\right) = \mu_{\mathbf{z}_i \mid \mathbf{x}_i} \quad \text{ and } \quad \mathsf{E}_{\mathbf{z}_i \mid \mathbf{x}_i}\left(z_i^2\right) = \mu_{\mathbf{z}_i \mid \mathbf{x}_i}^2 \ + \ \sigma_{\mathbf{z}_i \mid \mathbf{x}_i}^2$$

M-step:

$$\begin{split} \lambda_{j}^{\text{Gauss}} &= \left(\frac{1}{N}\sum_{i=1}^{N} x_{ij} \; \mathsf{E}_{\mathbf{z}_{i}|\mathbf{x}_{i}}\left(z_{i}\right) + \frac{1}{N}\frac{\mu_{\lambda} \; \Psi_{ij}^{\text{old}}}{\sigma_{\lambda}^{2}}\right) \left(\frac{1}{N}\sum_{i=1}^{N} \mathsf{E}_{\mathbf{z}_{i}|\mathbf{x}_{i}}\left(z_{i}^{2}\right) + \frac{1}{N}\frac{\Psi_{ij}^{\text{old}}}{\sigma_{\lambda}^{2}}\right)^{-1} \\ \lambda_{j}^{\text{new}} &= \begin{cases} \lambda_{j}^{\text{Gauss}} & \text{for } \lambda_{j}^{\text{Gauss}} > 0 \\ 0 & \text{for } \lambda_{j}^{\text{Gauss}} \leq 0 \end{cases}, \\ \Psi_{ij}^{\text{new}} &= \left[\operatorname{diagvect}\left(\frac{1}{N}\sum_{i=1}^{N} x_{i}x_{i}^{T}\right)\right]_{j} - \lambda_{j}^{\text{new}}\left[\frac{1}{N}\sum_{i=1}^{N} \mathsf{E}_{\mathbf{z}_{i}|\mathbf{x}_{i}}\left(z_{i}\right)x_{i}\right]_{j} + \frac{1}{N}\frac{\Psi_{ij}^{\text{old}}}{\sigma_{\lambda}^{2}}\lambda_{j}^{\text{new}}\left(\mu_{\lambda} - \lambda_{j}^{\text{new}}\right) \end{split}$$

◆ロト ◆昼 → ◆臣 → ◆臣 → ○ ◆ ◎ ◆

MA-plot FA vs. FARMS

FARMS



The MA plot shows log fold change as a function of mean log expression level.

FARMS

FARMS: Filtering by signal variance



(ロト (個) (目) (目) (日) (の)

FARMS: z-posterior

FARMS

Variance of $z \mid x$

Model

$$x = \mathbf{\lambda} z + \mathbf{\epsilon}$$

and Gaussian z-prior $\mathcal{N}(0,1)$ results in the z-posterior $p(z \mid x)$:

$$z \mid x \sim \mathcal{N} \left(\mu_{z \mid x} , \sigma_{z \mid x}^{2} \right)$$

$$\mu_{z \mid x} = (x)^{T} \Psi^{-1} \lambda \left(1 + \lambda^{T} \Psi^{-1} \lambda \right)^{-1}$$

$$\sigma_{z \mid x}^{2} = \left(1 + \lambda^{T} \Psi^{-1} \lambda \right)^{-1}$$

(ロ) (四) (日) (日) (日) (日) (日) (日)

FARMS: The I/NI call

FARMS

The variance of z is decomposed into a signal and a noise part:

$$1 = \operatorname{var}(z) = \frac{1}{N} \sum_{i=1}^{N} \operatorname{E}_{z_{i}|x_{i}}(z_{i}^{2}) = \frac{1}{N} \sum_{i=1}^{N} \left(\mu_{z_{i}|x_{i}}^{2} + \sigma_{z_{i}|x_{i}}^{2} \right)$$
$$\frac{1}{N} \sum_{i=1}^{N} \sigma_{z_{i}|x_{i}}^{2} = 1 - \frac{1}{N} \sum_{i=1}^{N} \mu_{z_{i}|x_{i}}^{2}$$
$$\sigma_{z|x}^{2} = 1 - \frac{1}{N} \sum_{i=1}^{N} \mu_{z_{i}|x_{i}}^{2} = \left(1 + \boldsymbol{\lambda}^{T} \boldsymbol{\Psi}^{-1} \boldsymbol{\lambda} \right)^{-1}$$

 $\sigma_{z|x}^2$ is called the "Informative/NonInformative (I/NI) call" and is one minus the signal variance. We see that large λ (going with low noise Ψ) leads to low variance of $z \mid x$ which means a precise conditional z.

・ロト ・聞 ト ・ヨト ・ヨー うへぐ

FARMS

FARMS: Independent I/NI calls filtering

Independent filtering increases detection power for high-throughput experiments

Richard Bourgon^a, Robert Gentleman^b, and Wolfgang Huber^{c,1}

"European Bioinformatics Institute, Cambridge CB10 ISD, United Kingdom; ^bGenentech, Inc., 1 DNA Way, South San Francisco, CA 94080-4990; and 'European Molecular Biology Laboratory, 69117 Heidelberg, Germany

Edited by Stephen E. Fienberg, Carnegie Mellon University, Pittsburgh, PA, and approved March 22, 2010 (received for review December 3, 2009)

- For permutation invariant test statistics and for the *t*-test statistic *T* (only for Gaussian *z*-prior), the I/NI call filter applied to null hypotheses is independent of the statistic
- This guarantees type I error rate control if first filtering by I/NI calls, then using these statistics, and finally applying correction for multiple testing.
- http://www.bioinf.jku.at/software/cnfarms/proof_ini.pdf

FARMS: I/NI calls distribution



Bimodal distribution

- Enforced by the parameter prior
- Modes clearly separated (insensitive for filtering threshold)
- Works for unbalanced data (few samples contain a signal) in contrast to variance filtering (Bourgon et al. (2010))
- Works for few genes with a signal

FARMS

A pipeline for gene expression analysis



FARMS

Receiver Operator Characteristics (ROC)

Affycomp II / GoldenSpike Benchmark (AUC - area under the curve):								
	INTENSITY	FARMS	RMA	GCRMA	MAS 5.0	MBEI		
HGU133	Low	0.94	0.51	0.62	0.07	0.21		
	Med	0.99	0.91	0.94	0.00	0.43		
	High	1.00	0.64	0.59	0.00	0.16		
	Mean	0.95	0.60	0.69	0.05	0.26		
HGU95	Low	0.91	0.57	0.45	0.09	-		
	Med	1.00	0.91	0.91	0.00	-		
	High	0.98	0.96	0.92	0.00	-		
	Mean	0.93	0.65	0.57	0.06	-		
GoldenSpike		0.85	0.76	0.78	0.28	0.39		

A CC 1 (11 - 0 . . . ١

Computational costs for processing 60 arrays

•	-	-	-		
		FARMS	RMA	MAS 5.0	MBEI
C	OMPUTATIONAL TIME [S]	92	384	851	591

Results I/NI call

FARMS

- \bullet Leads on average to 84 (±1.5)% exclusion rate
 - Applied on 30 real life studies
 - A/P calls excluded only 33 $(\pm 1)\%$
- Validation was carried out on spiked-in data:

Exclusion rate on spiked-in data sets:

	INFORMATIVE	Non-informative	Exclusion rate	Detected Spiked-ins	Detected Pseudo Spiked-ins
HGU133A	81	22219	99.63%	42/42	28/28*
HGU95_V2	56	12570	99.56%	14/14	5/5**
HU. GENE 1.0 ST	40	19,753	99.80%	15/15***	-

*McGee et al. 2006; **Wolfinger and Chu 2002; Cope et al. 2004; ***long spiked-in fragments

Copy number variants



- are deletions, duplications, inversions or inserations of chromosomal segments
- are a major source of variation between individual humans
- are an underlying factor in human evolution and in many diseases
- form at a faster rate than other types of mutation

Rare CNV events

Sparse data

CNV data is sparse with an kurtosis larger than 30 \rightarrow change the model assumption to a Laplacian distributed hidden variable z.





Close up Gauss vs. Laplace



Laplacian cn.FARMS

Data likelihood

$$p(\{x\} | \boldsymbol{\lambda}, \boldsymbol{\Psi}) = \int p(\{x\} | z, \boldsymbol{\lambda}, \boldsymbol{\Psi}) p(z) dz$$

Laplacian cn.FARMS

Data likelihood

$$p(\{x\} \mid \boldsymbol{\lambda}, \boldsymbol{\Psi}) = \int p(\{x\} \mid z, \boldsymbol{\lambda}, \boldsymbol{\Psi}) p(z) dz$$

Problem

• The **likelihood is analytically intractable** for the non-Gaussian prior

Laplacian cn.FARMS

Data likelihood

$$p(\{x\} \mid \boldsymbol{\lambda}, \boldsymbol{\Psi}) = \int p(\{x\} \mid z, \boldsymbol{\lambda}, \boldsymbol{\Psi}) p(z) dz$$

Problem

• The **likelihood is analytically intractable** for the non-Gaussian prior

Solution

- Variational EM approach
- Based on a local Gaussian approximation to the mode

Benchmark data sets



- 30 male and 30 female CEU founders
 - Classification task: distinguish males from females by their copy number on the X chromosome
- Evaluation on:
 - Single-locus / multi-loci classification (window mode)
 - Multi-loci summarization with
 - cn.FARMS
 - Median locus for dChip and CRMA_v2

ROC-Curve (SNP 6.0 arrays)



TPR / FPR

True positive rate (TPR) = TP/(TP+FN)False positive rate (FPR) = FP/(FP+TN)

Results cn.FARMS

		AFFYMETE	rix Mapping2 <u>5</u>	oK_NSP	Affymetrix SNP 6.0			
Loci	Criteria	cn.FARMS	CRMA_v2	dChip	cn.FARMS	CRMA_v2	dChip	
1	AUC	0.9852	0.9820	0.9819	0.9838	0.9807	0.9721	
	FP	8472	9106	9018	56145	68593	77438	
	P-VALUE	-	1.8e-65	3.1e-26	-	1e-1160	1e-6049	
2	AUC	0.9983	0.9974	0.9969	0.9983	0.9963	0.9894	
	FP	1375	1449	1611	9777	11705	18039	
	P-VALUE	-	2.7e-4	2.5e-12	-	1e-317	1e-3713	
3	AUC	0.9998	0.9995	0.9992	0.9998	0.9990	0.9953	
	FP	240	366	440	1573	3462	6625	
	P-VALUE	-	2.6e-38	7.2e-58	-	1e-896	1e-3455	
	AUC	1.000	0.9999	0.9998	0.9999	0.9995	0.9976	
Tabla : AUC values at the sex classification task for 590 Hap Map 30 EU founders								
based on t	the-X _{LG} hro	omosome o	-	1e-594	1e-2013			

CNV detection benchmark

- "The International HapMap Project" phase 2 data set with Affymetrix SNP 6.0 arrays
 - $\bullet\,$ Goal is to identify true rare CNV regions with a low FDR
 - "True CNV regions" are those regions which were detected and verified by different bio-technologies
 - NimbleGen tiling arrays, Agilent CGH arrays, Illumina Infinium genotyping (Human660W)
 - 2,515 true CNV regions as reference
- CNV calling criteria:
 - I/NI call for cn.FARMS
 - $\bullet\,$ Variance of the raw copy numbers on the samples for dChip and CRMA_v2

CNV detection plot



CNV detection on HapMap (multi-loci 3)

Chromosome 8



Whole genome



Precision / Recall

 $\begin{aligned} \text{Recall} &= \text{TP}/(\text{TP}+\text{FN}) \\ \text{Precision} &= \text{TP}/(\text{TP}+\text{FP}) = 1 \text{ - FDR} \end{aligned}$

200

CNV detection on HapMap (multi-loci 5)

Chromosome 8



Whole genome



Precision / Recall

 $\begin{aligned} \text{Recall} &= \text{TP}/(\text{TP}+\text{FN}) \\ \text{Precision} &= \text{TP}/(\text{TP}+\text{FP}) = 1 \text{ - FDR} \end{aligned}$

900

ICGC copy number data sets

- Glioblastoma multiforme data sets
 - 167 Agilent 415K CGH arrays from Harvard
 - 262 Agilent 244A CGH arrays from Harvard
 - 461 Agilent 244A CGH arrays from MSKCC
 - 533 Affymetrix SNP 6.0 arrays from Broad
 - 432 Illumina HumanHap 550 from Stanford
- CN data for SNP 6.0 and HumanHap 550 were not available
- 167 matched arrays HMS 415K and MSKCC 244A remain
Merged raw data (Chromosome 1)



コント (日) (日) (日) (日) (日) (日)

Prior weight 2.0



◆ロ → ◆昼 → ◆ 臣 → ◆ 臣 → ⑦ � ⑦

Biclustering applications

Definition: Biclustering simultaneously organizes a data matrix into subsets of rows and columns in which the entities of each row subset are similar to each other on the column subset and vice versa.

- $\bullet~{\sf Gene}~{\sf expression}$ $\Rightarrow~{\sf columns}$ tissues, rows genes
 - compounds that trigger the same pathway
 - tightly co-expressed gene sets in subgroups of cancer, e.g. patients with bad treatment outcome
- $\bullet~{\sf Bioassays} \Rightarrow~{\sf columns}$ compounds, rows bioassay activity
 - compounds that are active on similar targets
- $\bullet~Structural~fingerprints \Rightarrow~columns~compounds,~rows~chemical~fingerprints$
 - compounds that share a chemical substructure

Biclustering: The idea



(ロ) (四) (三) (三) (三) (三) (○)

FABIA: The model I

factor z

loading matrix

observations x

poise ϵ



◆□ → ◆■ → ◆臣 → ◆臣 → ○ ● ● ● ● ●

FABIA: The model II

$$x = \mathbf{\lambda}z + \mathbf{\epsilon} = \sum_{i=1}^{p} \lambda_i z_i^T + \mathbf{\epsilon}$$

- x are the observations
- $\bullet~\lambda$ is the matrix of factor loadings
- $z = (z_1, \ldots, z_p)^T$ is the factor matrix
- *p* number of biclusters
- $\lambda_i \in \mathbb{R}^n$ is the sparse prototype vector of the *i*-th bicluster
- $z_i \in \mathbb{R}^l$ is the sparse vector of factors of the *i*-th bicluster
- $\epsilon \in \mathbb{R}^{n \times l}$ is independent additive noise $\mathscr{N}(0, \Psi)$ -distributed

FABIA: Bayes framework

Loading prior

- Sparseness on the loadings
- Laplace prior

•
$$p(\mathbf{\lambda}_i) = \left(\frac{1}{\sqrt{2}}\right)^n \prod_{j=1}^n e^{-\sqrt{2}|\lambda_{jj}|}$$

FABIA: Bayes framework

Loading prior

- Sparseness on the loadings
- Laplace prior

•
$$p(\mathbf{\lambda}_i) = \left(\frac{1}{\sqrt{2}}\right)^n \prod_{j=1}^n e^{-\sqrt{2}|\lambda_{ji}|}$$

Factor prior

- Sparseness on the factor
- Laplace prior

•
$$p(z) = \left(\frac{1}{\sqrt{2}}\right)^p \prod_{i=1}^p e^{-\sqrt{2}|z_i|}$$

FABIA: Bayes framework

Loading prior

- Sparseness on the loadings
- Laplace prior

•
$$p(\mathbf{\lambda}_i) = \left(\frac{1}{\sqrt{2}}\right)^n \prod_{j=1}^n e^{-\sqrt{2}|\lambda_{ji}|}$$

Factor prior

- Sparseness on the factor
- Laplace prior

•
$$p(z) = \left(\frac{1}{\sqrt{2}}\right)^p \prod_{i=1}^p e^{-\sqrt{2}|z_i|}$$

Problem

Laplace prior on factors leads to intractable likelihood:

$$p(x \mid \boldsymbol{\lambda}, \boldsymbol{\Psi}) = \int p(x \mid z, \boldsymbol{\lambda}, \boldsymbol{\Psi}) p(z) dz$$

Solution: Prior on factors is replaced by maximum of a Gaussian function family \Rightarrow variational approach

$$p(z) \approx \underset{\xi}{\operatorname{arg\,max}} p(z|\xi)$$

FABIA: Variational EM updates

E-step:

$$\begin{split} & \mathbb{E}(\tilde{z}_{j} \mid x_{j}) \ = \ \left(\boldsymbol{\lambda}^{\mathcal{T}} \ \boldsymbol{\Psi}^{-1} \ \boldsymbol{\lambda} \ + \ \boldsymbol{\Xi}_{j}^{-1}\right)^{-1} \ \boldsymbol{\lambda}^{\mathcal{T}} \ \boldsymbol{\Psi}^{-1} \ x_{j} \quad \text{and} \\ & \mathbb{E}(\tilde{z}_{j} \ \tilde{z}_{j}^{\mathcal{T}} \mid x_{j}) \ = \ \left(\boldsymbol{\lambda}^{\mathcal{T}} \ \boldsymbol{\Psi}^{-1} \ \boldsymbol{\lambda} \ + \ \boldsymbol{\Xi}_{j}^{-1}\right)^{-1} \ + \ \mathbb{E}(\tilde{z}_{j} \mid x_{j}) \ \mathbb{E}(\tilde{z}_{j} \mid x_{j})^{\mathcal{T}} \quad \text{where} \\ & \mathbf{\Xi}_{j} \ = \ \text{diag}(\text{diagvect}\left(\sqrt{\mathbb{E}(\tilde{z}_{j} \ \tilde{z}_{j}^{\mathcal{T}} \mid x_{j})}\right)) \text{ is the update for the variational parameter.} \end{split}$$

FABIA: Variational EM updates

E-step:

$$\begin{split} & \mathbb{E}(\tilde{z}_{j} \mid x_{j}) \ = \ \left(\boldsymbol{\lambda}^{\mathcal{T}} \ \boldsymbol{\Psi}^{-1} \ \boldsymbol{\lambda} \ + \ \boldsymbol{\Xi}_{j}^{-1}\right)^{-1} \ \boldsymbol{\lambda}^{\mathcal{T}} \ \boldsymbol{\Psi}^{-1} \ x_{j} \quad \text{and} \\ & \mathbb{E}(\tilde{z}_{j} \ \tilde{z}_{j}^{\mathcal{T}} \mid x_{j}) \ = \ \left(\boldsymbol{\lambda}^{\mathcal{T}} \ \boldsymbol{\Psi}^{-1} \ \boldsymbol{\lambda} \ + \ \boldsymbol{\Xi}_{j}^{-1}\right)^{-1} \ + \ \mathbb{E}(\tilde{z}_{j} \mid x_{j}) \ \mathbb{E}(\tilde{z}_{j} \mid x_{j})^{\mathcal{T}} \quad \text{where} \\ & \mathbf{\Xi}_{j} \ = \ \text{diag}(\text{diagvect}\left(\sqrt{\mathbb{E}(\tilde{z}_{j} \ \tilde{z}_{j}^{\mathcal{T}} \mid x_{j})}\right)) \text{ is the update for the variational parameter.} \end{split}$$

M-step:

$$\begin{split} \lambda^{\mathrm{new}} &= \frac{\frac{1}{l} \sum_{j=1}^{l} x_j \ \mathbb{E}(\tilde{z}_j \mid x_j)^{\mathcal{T}} - \frac{\alpha}{l} \ \Psi \ \mathrm{sign}(\lambda)}{\frac{1}{l} \sum_{j=1}^{l} \mathbb{E}(\tilde{z}_j \ \tilde{z}_j^{\mathcal{T}} \mid x_j)} \\ \Psi^{\mathrm{new}} &= \Psi^{\mathrm{EM}} + \mathrm{diag}(\mathrm{diagvect}\Big(\frac{\alpha}{l} \ \Psi \ \mathrm{sign}(\lambda)(\lambda^{\mathrm{new}})^{\mathcal{T}}\Big)\Big) \ , \quad \text{where} \\ \Psi^{\mathrm{EM}} &= \mathrm{diag}(\mathrm{diagvect}\Big(\frac{1}{l} \sum_{j=1}^{l} x_j x_j^{\mathcal{T}} - \lambda^{\mathrm{new}} \frac{1}{l} \sum_{j=1}^{l} \mathbb{E}\left(\tilde{z}_j \mid x_j\right) \ x_j^{\mathcal{T}}\Big) \Big). \end{split}$$

 α controls the degree of sparseness (parameter of the Laplacian prior)

Biclustering of copy number variants I



◆ロト ◆聞 ▶ ◆臣 ▶ ◆臣 ▶ ◆ 国 ● のへの

Biclustering of copy number variants II



< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Biclustering of bioassays and compounds

Matrix plot (dim 270,000 × 4,000)



Bioassay data details

- Data source: ChEMBL
- # of assays: ca. 4,000
- # of compounds: ca. 270,000
- Sparseness: ca. 1:2,000

Biclustering of bioassays and compounds

Matrix plot - close up



Bioassay data details

- Data source: ChEMBL
- # of assays: ca. 4,000
- # of compounds: ca. 270,000
- Sparseness: ca. 1:2,000

Biclustering of bioassays and compounds



Compounds

Compounds of the bioassay bicluster



◆ロ → ◆聞 → ◆臣 → ◆臣 → ○ ● ○ ○ ○ ○

Biclustering of fingerprints and compounds

Matrix plot (dim $16e+6 \times 1e+6$)



Bioassay data details

- Data source: ChEMBL
- # of fingerprints: ca. 16,000,000
- # of compounds: ca. 1,000,000
- Sparseness: ca. 1:150,000

Bicluster of fingerprints and compounds I



・ロト ・母 ト ・目 ト ・日 ・ うへぐ

Bicluster of fingerprints and compounds I



Compounds of the fingerprint bicluster













Compounds of the fingerprint bicluster



• All compounds of this bicluster show kinase bioactivity (urokinase-type plasminogen activator)

Biclustering for recommender systems

Matrix plot (dim many x many)

Recommender

- Data source: Zalando
- # of articles: many
- # of cookies: many
- Sparseness: ca. 1:20,000

Overview 12 random selected bicluster



200

Top 36 articles of a bicluster



- emerging machine learning technique
- $\bullet\,$ multiple levels of sparse representations $\Rightarrow\,$ higher levels representing more abstract concepts
- Google and facebook now apply deep learning for object recognition, image and information retrieval
- Google recently acquired the deep learning start-up DeepMind for \$500M, winning bidding against facebook
- Nature and The New York Times covered deep learning with serveral articles (two front-page articles)

Networks



◆□ > ◆圖 > ◆国 > ◆国 > → 国 → のへで

Sparsity by Linear-Rectified Units



Sparsity by Dropout I



- randomly set units to zero acitvation
- no derivatives

Sparsity by Dropout II



Model: Rectified Factor Network



Rectified Factor Network: Newton updates

Newton-step:

$$\begin{split} \boldsymbol{\mu}_{h|\boldsymbol{\nu}} &= \boldsymbol{W}^{T} \left(\boldsymbol{W} \, \boldsymbol{W}^{T} \, + \, \boldsymbol{\Psi} \right)^{-1} \, \boldsymbol{\nu} \, , \\ \boldsymbol{\Sigma}_{h|\boldsymbol{\nu}} &= \boldsymbol{I} - \boldsymbol{W}^{T} \left(\boldsymbol{W} \, \boldsymbol{W}^{T} \, + \, \boldsymbol{\Psi} \right)^{-1} \boldsymbol{W} \, , \\ \boldsymbol{E}_{h_{i}|\boldsymbol{\nu}_{i}} \left(h_{i} \right) &= \boldsymbol{\mu}_{h_{i}|\boldsymbol{\nu}_{i}} \\ \boldsymbol{E}_{h_{i}|\boldsymbol{\nu}_{i}} \left(h_{i} \, h_{i}^{T} \right) &= \boldsymbol{\mu}_{h_{i}|\boldsymbol{\nu}_{i}} \, \boldsymbol{\mu}_{h_{i}|\boldsymbol{\nu}_{i}}^{T} + \boldsymbol{\Sigma}_{h_{i}|\boldsymbol{\nu}_{i}} \\ \boldsymbol{U} &= \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{\nu}_{i} \boldsymbol{E}_{h_{i}|\boldsymbol{\nu}_{i}}^{T} \left(h_{i} \right) \quad \text{Hebb rule: input - hidden} \\ \boldsymbol{S} &= \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{E}_{h_{i}|\boldsymbol{\nu}_{i}} \left(h_{i} \, h_{i}^{T} \right) \end{split}$$

Gradients:

 $\Delta W = U S^{-1} - W \text{ Hebb rule: decorrelates factors}$ $\Delta \Psi_{ii} = \left[C - U W^{T} - W U^{T} + W S W^{T} \right]_{ii} - \Psi_{ii}$

MNIST data set

MNIST 28 x 28 pixel



MNIST

- Data source: Yann LeCun
- # of input pixel: 756
- # of samples: 70,000
- Classify digits

・ロト ・聞 ・ ・ ヨ ・ ・ ヨ ・ のへで

Benchmark data



(a) rot, bg-rand, bg-img, bg-img-rot

Samples form the various image classification problems. (a): harder variations on the MNIST digit classification problems. (b): artificial binary classification problems.

200

⁽b) rect, rect-img, convex

Receptive fields



Various filters learnt from 1024 hidden units RFN on benchmark data set.

Motivation

IS cn.F/

Biclusterin

Results

Dataset		SVM_{rbf}	DBN ₁	DBN ₃	SAE ₃	SDAE ₃	RFN	
MNIST	50k-10k-10k	$1.40{\scriptstyle \pm 0.23}$	1.21 ± 0.21	$1.24{\scriptstyle \pm 0.22}$	$1.40{\scriptstyle \pm 0.23}$	1.28 ± 0.22	$1.27{\scriptstyle\pm0.22}$	(1)
basic	10k-2k-50k	$3.03{\scriptstyle \pm 0.15}$	3.94 ± 0.17	$3.11{\scriptstyle \pm 0.15}$	3.46 ± 0.16	2.84 ± 0.15	$2.66{\scriptstyle \pm 0.14}$	(1)
bg-rand	10k-2k-50k	14.58 ± 0.31	$9.80{\scriptstyle \pm 0.26}$	$6.73{\scriptstyle \pm 0.22}$	11.28 ± 0.28	10.30 ± 0.27	$7.94{\scriptstyle \pm 0.24}$	(3)
bg-img	10k-2k-50k	22.61 ± 0.37	$16.15{\scriptstyle \pm 0.32}$	$16.31{\scriptstyle \pm 0.32}$	23.00±0.37	$16.68{\scriptstyle \pm 0.33}$	$16.52{\scriptstyle \pm 0.32}$	(1)
rect	1k-0.2k-50k	2.15 ± 0.13	$4.71{\scriptstyle \pm 0.19}$	$2.60{\scriptstyle \pm 0.14}$	2.41 ± 0.13	1.99 ± 0.12	0.63 ± 0.06	(1)
rect-img	10k-2k-50k	24.04 ± 0.37	$23.69{\scriptstyle \pm 0.37}$	22.50 ± 0.37	24.05±0.37	21.59 ± 0.36	$20.77{\scriptstyle\pm0.36}$	(1)
convex	10k-2k-50k	$19.13{\scriptstyle \pm 0.34}$	$19.92{\scriptstyle \pm 0.35}$	$18.63{\scriptstyle\pm0.34}$	$18.41{\scriptstyle \pm 0.34}$	$19.06{\scriptstyle \pm 0.34}$	$16.41{\scriptstyle\pm0.32}$	(1)

Test error rate on all considered classification problems is reported together with a 95% confidence interval.
Predicting Drug-Target Interactions

Target prediction



AUC - Area under the ROC curve

- Data source: ChemBL
- Compounds: 698,425
- Targets: 1,230
- Descriptors: 43,340
- Hidden units: 16,384
- Parameters: 422,232,064
- Computation: Nvidia Tesla K40 with 2,880 CUDA GPU cores

- Sparse coding can reliably identify interesting projection in the data
- Sparse coding can be used for biclustering of high-dimensional data
- Sparse coding in drug design can help in selecting compounds with strong on-target effects and thereby helps to impute missing measurements
- Rectified linear units in combination with dropout lead to sparse representations of the data
- Rectified Factor Networks outperform all existing unsupervised deep learning methds and can be used for various problems

Open source software



- FARMS, cn.FARMS and FABIA are publicly available as Bioconductor R packages
- Software homepages:
 - http://www.bioinf.jku.at/software/farms/farms.html
 - http://www.bioinf.jku.at/software/cnfarms/cnfarms.html
 - http://www.bioinf.jku.at/software/fabia/fabia.html