VI - simple cell receptive fields



combine Simple cells sum LGN inputs



The "standard model" of VI



Sparse coding model



$$P(\mathbf{a}|\mathbf{I}; \mathbf{\Phi}) \propto P(\mathbf{I}|\mathbf{a}; \mathbf{\Phi}) P(\mathbf{a})$$
$$\hat{\mathbf{a}} = \arg\min_{\mathbf{a}} |\mathbf{I} - \mathbf{\Phi} \mathbf{a}|^2 + \lambda \sum_{i} C(a_i)$$
$$Not \quad \hat{a}_i = g(\sum_{x,y} \phi_i(x, y) I(x, y))$$

- Cortex is *not* doing PCA
 - PCA only captures pairwise correlations in data.
 - PCA assumes Gaussianity and orthogonal components.
 - Principal components do not resemble simple-cell RFs.
- Cortex is *not* doing compression
 - expands dimensionality of representation.
 - goal is to *interpret* image data, not to simply encode or compress it.
- Cortex seems to be something akin to sparse coding
 ...and probably much more.

Autoencoder networks



PCA (Principal Components Analysis)



 $\langle x_1 \, x_2 \rangle = c_{12} \qquad y_1 = \mathbf{e}_1 \cdot \mathbf{x} \qquad \langle y_1 \, y_2 \rangle = \langle y_1 \rangle \, \langle y_2 \rangle$ $\neq 0 \qquad y_2 = \mathbf{e}_2 \cdot \mathbf{x} \qquad = 0$

$$\mathbf{E} = \begin{bmatrix} | & | \\ \mathbf{e}_1 & \mathbf{e}_2 \\ | & | \end{bmatrix} \qquad \mathbf{e}_1 \cdot \mathbf{e}_2 = 0$$
$$|\mathbf{e}_1| = |\mathbf{e}_2| = 1$$

PCA (Principal Components Analysis) *a*.

b





 $\Delta W_{AB} \propto \langle AB \rangle$

Linear Hebbian learning



 $\dot{w}_i \propto \langle y x_i \rangle$ $y = \sum_{j} w_j x_j$ $\dot{w}_i \propto \left\langle \sum_i w_j x_j x_i \right\rangle$ $= \sum w_j \langle x_j x_i \rangle$

 $\dot{\mathbf{w}} \propto \mathbf{C} \mathbf{w}$ $C_{ij} = \langle x_i \, x_j \rangle$

$\mathbf{C} = \mathbf{E} \mathbf{\Lambda} \mathbf{E}^T$







Principal components of natural image patches (8 x 8 pixels)

$$\mathbf{W} = \mathbf{E}^T$$



- Not localized
- Not oriented

PCA is incapable of learning about localized, oriented structure in images.

I/f noise

(what the world looks like if all you care about are pairwise correlations)



Higher-order image statistics

phase alignment orientation

motion



Bottleneck may also be in the form of limited capacity units. Optimal strategy in this case is to whiten.



Efficient coding model of retina (Karklin & Simoncelli 2012)





Sparse codes impose a different type of bottleneck by limiting the number of active units





Evidence for grandmother cells?

(Quiroga, Reddy, Kreiman, Koch & Fried, *Nature* 2005)



Evidence for grandmother cells?

(Quiroga, Reddy, Kreiman, Koch & Fried, Nature 2005)



Evidence for grandmother cells?

(Quiroga, Reddy, Kreiman, Koch & Fried, Nature 2005)



Winner-take-all learning



Learning rule: $\Delta w_{ij} = \eta \, y_i \, (x_j - w_{ij})$

Winner-take-all learning

before learning

after learning





Sparse, distributed representation



Biological Cybernetics

Forming sparse representations by local anti-Hebbian learning

P. Földiák

Physiological Laboratory, University of Cambridge, Downing Street, Cambridge CB2 3EG, United Kingdom

$$\frac{dy_i^*}{dt} = f\left(\sum_{j=1}^m q_{ij}x_j + \sum_{j=1}^n w_{ij}y_j^* - t_i\right) - y_i^*$$



anti-Hebbian rule-

 $\Delta t_i = \gamma(y_i - p) \; .$

$$\Delta w_{ij} = -\alpha (y_i y_j - p^2)$$

(if $i = j$ or $w_{ij} > 0$ then $w_{ij} := 0$)
Hebbian rule-
 $\Delta q_{ij} = \beta y_i (x_j - q_{ij})$
threshold modification-

Learning lines

Input patterns:



Learned weights:



PCA solution



Reconstructions



Problems

- How to deal with graded input signals? (i.e., real images)
- No objective function

Sparse coding model (Olshausen & Field 1996, 1997)



$$I_{I(x,y)} \bigoplus_{\phi_i(x,y)} \phi_i(x,y) = \sum_i a_i \phi_i(x,y) + \epsilon(x,y)$$

$$P(\mathbf{a}|\mathbf{I}; \Phi) \propto P(\mathbf{I}|\mathbf{a}; \Phi) P(\mathbf{a})$$

$$\hat{\mathbf{a}} = \arg \min_{\mathbf{a}} |\mathbf{I} - \Phi \mathbf{a}|^2 + \lambda \sum C(a_i)$$

i

$P(\mathbf{a}|\mathbf{I}; \mathbf{\Phi}) \propto P(\mathbf{I}|\mathbf{a}; \mathbf{\Phi}) P(\mathbf{a})$

Cost function

$$C(a_i) = \log(1 + a_i^2)$$





ai

Compute coefficients via gradient descent

$$\tau \dot{a}_i = -\frac{dE}{da_i}$$
$$= b_i - \sum_{j \neq i} G_{ij} a_j - f_\lambda(a_i)$$

Where

$$b_{i} = \sum_{x,y} \phi_{i}(x,y) I(x,y)$$
$$G_{ij} = \sum_{x,y} \phi_{i}(x,y) \phi_{j}(x,y)$$
$$f_{\lambda}(a_{i}) = a_{i} + \lambda C'(a_{i})$$

Alternative formulation (the Hopfield trick)

Let

$$u_{i} = f_{\lambda}(a_{i}), \text{ or } a_{i} = f_{\lambda}^{-1}(u_{i}) \equiv g(u_{i})$$
$$\tau \dot{u}_{i} = -\frac{dE}{da_{i}}$$
$$= b_{i} - \sum_{j \neq i} G_{ij} a_{j} - u_{i}$$

Thus

$$\tau \dot{u}_i + u_i = b_i - \sum_{j \neq i} G_{ij} a_j$$
$$a_i = g(u_i)$$

Relation between the thresholding function g and cost function C



Sparse inference via lateral inhibition and thresholding

(Rozell, Johnson, Baraniuk & Olshausen, 2008)



$$\tau \dot{u}_i + u_i = b_i - \sum_{j \neq i} G_{ij} a_j$$
$$a_i = g(u_i)$$

$$b_{i} = \sum_{\mathbf{x}} \phi_{i}(\mathbf{x}) I(\mathbf{x})$$
$$G_{ij} = \sum_{\mathbf{x}} \phi_{i}(\mathbf{x}) \phi_{j}(\mathbf{x})$$

Energy function

$$E = \frac{1}{2} |\mathbf{I} - \Phi \mathbf{a}|^2 + \lambda \sum_{i} C(a_i)$$

Learning rule

$$\Delta \phi_i = -\eta \frac{\partial E}{\partial \phi_i} = [\mathbf{I} - \Phi \, \hat{\mathbf{a}}] \, \hat{a}_i$$

$\Phi_i(x, y)$ learned from natural images (200, 12x12 pixels)





linear projection sparsified coefficients image reconstruction

					-									1													
								4		/																	
		1	8																								
					-																111						
													11														
			+	1																			10				
			+												- 55		0	111								\rightarrow	
			+																							\rightarrow	
			+																								
																					-	1				\rightarrow	
			+																						-	\rightarrow	
	~		+					1		-	1	9		•	1					1							
			+						_															X			
		-	+																						-	\rightarrow	
			+				11											24.					_			\rightarrow	
	-	-	+																								
10			+								- 4								-						_	\rightarrow	
			+				~					1															
		_	+									-	-													-	
	_	-	+																			/			1	-	
			+															1							_		
			+											630											- 4		
			+													1.0									_	\rightarrow	
		_	+													1										\rightarrow	
			$ \rightarrow $						1				1														_
1																											
						1																					
																	-										
											1	-															
							-							-					NW S								

Deep convnets are easily fooled by imperceptible perturbations (adversarial examples)



Szegedy et al. (2013)

Sparse inference protects against adversarial attack (Paiton, Frye, Lundquist, Bowen, Zarcone & Olshausen 2020)



Dylan Paiton



Joel Bowen

iso-response contours



linear projection

sparsified

Relationship between Sparse Coding and ICA

• No noise

$$\mathbf{x} = \mathbf{A}\mathbf{s}$$

• Invertible **A** matrix

$$\mathbf{s} = \mathbf{A}^{-1}\mathbf{x}$$

Thus
$$p(\mathbf{x}|\mathbf{s}) = \delta(\mathbf{x} - \mathbf{A}\mathbf{s})$$

$$p(\mathbf{x}) = \int \delta(\mathbf{x} - \mathbf{A} \mathbf{s}) p_s(\mathbf{s}) d\mathbf{s}$$
$$= p_s(\mathbf{A}^{-1}\mathbf{x}) / |\det \mathbf{A}|$$

$$\log p(\mathbf{x}) = -\sum_{i} C(s_i) - \log \det \mathbf{A}$$

$$\Delta \mathbf{A} \propto \frac{\partial}{\partial \mathbf{A}} \langle \log p(\mathbf{x}) \rangle \qquad \begin{array}{l} \text{Its the ICA} \\ \text{learning rule!} \\ = \frac{\partial}{\partial \mathbf{A}} \left[-\sum_{i} C(s_i) - \log \det \mathbf{A} \right] \\ = \left[\langle [\mathbf{A}^T]^{-1} \mathbf{z}(\mathbf{s}) \mathbf{s}^T - [\mathbf{A}^T]^{-1} \rangle \right] \checkmark$$

Pre-multiplying by $\mathbf{A} \mathbf{A}^{\mathsf{T}}$ (natural gradient) yields:

$$\Delta \mathbf{A} \propto \langle \mathbf{A} \mathbf{z} \mathbf{s}^T - \mathbf{A} \rangle$$

= $\langle [\mathbf{x} - \mathbf{A}(\mathbf{s} - \mathbf{z})] \mathbf{s}^T - \mathbf{A} \rangle$

The "Independent Components" of Natural Scenes are Edge Filters

ANTHONY J. BELL,*† TERRENCE J. SEJNOWSKI* Received 16 July 1996; in revised form 9 April 1997





EEG

Fz	man Martin Marine Ma
Cz	man when when a second second
Pz	www.WWWWWWWWWWWW
Oz	www.ww
F3	man My Mary Mary Mary Market Market
F4	man have man and have a second the second se
C3	www.www.www.www.www.www.www.www.
C4	man Mar
T3	annor marked of the the the the the second and the
T4	Mayne Marine Ma
P3	www.www.www.www.www.www.www.
P4	www.Www.www.www.www.
Fpz	www.www.www.www.www.www.www.www.
EOG	man was warden and the second se

ICA

 $1 \dots 1$

- 3 white Marken marken marken and the second second
- 4 www.www.www.www.www.www.www.

- 7 manufamana Manumanana
- 8 MMM And Marken M Marken Mark
- 10
- 11 Martin Manus Martin Martin
- 13 and the second secon
- 14 + 1 sec.

theta wave alpha wave

EOG

line noise

ICA is not a general solution for finding independent components of data

Assumptions of the model:

I. linear superposition: $\mathbf{x} = \mathbf{As}$

2. shape of the prior over each of the components: $p(s_i) \propto e^{-C(s_i)}$

3. factorial prior over the entire set of components:

$$p_s(\mathbf{s}) \propto \prod_i p_{s_i}$$