

NEURAL COMPUTATION VS 265

CONTENTS

1.1	Sparse Distributed Coding	2
1.2	Foldiak Paper - Sparse Coding	5
1.3	Comprehensive Review	8
1.3.1	Unsupervised Learning	8
1.4	Lab 4 & LCA Handout	12
1.5	HKP 9.4 - Feature Mapping	14
1.6	Locally Linear Embedding	15
1.7	Recurrent Neural Networks	18
1.8	Hopfield Networks Handout	21
1.9	Mixture of Gaussians and EM Algorithm	23
1.10	Boltzmann Machines	25
1.10.1	Lecture Slides	25
1.10.2	HKP Chapter 7.1	25
1.10.3	Boltzmann Machines - AIFH	26
1.11	Independent Component Analysis	27
1.11.1	ICA - Andrew Ng - CS 229	29

Sparse Distributed Coding

Table of Contents Local

Written by Brandon McKinzie

- VI simple-cell receptive fields are localized, oriented, and bandpass.
- PCA is really bad for such situations.
- To detect sharp edges in images, need high frequency and in-phase combinations.
- Higher-order image statistics:
 - phase alignment
 - orientation
 - motion
- want to move beyond pairwise correlations.
- WTA is too greedy, want more distributed strategy.
- Idea: **Projection pursuit**.
 - Look for low-dimensional projections that are as non-Gaussian as possible.
 - Projections tend to result in Gaussian distributions by the C.L.T.
 - Want to explore projections onto a weight vector until find something Non-Gaussian. Why? Because such a distribution could not have happened by accident.
- **Gabor-filter** response histograms are highly non-Gaussian.
- (Lab-related) Paper on *Forming sparse representations by local anti-Hebbian learning*.
 - Each neuron takes weighted input sum, as well as getting lateral inhibition by neighbors, but where the lateral weights are all negative. Put all through f , some sigmoidal non-linearity. "leaky integrator"
 - Want population-sparsity, so need neurons decorrelated. Have three learning rules: anti-Hebbian, Hebbian, and threshold modification.

- Threshold modification resembles homeostasis.

$$\Delta t_i = \gamma(y_i - p) \quad (1)$$

which is essentially SGD. Think about average behavior, as it relates to y_i output and p . p is a constant to be determined. Feedback loop. Adjusts spiking threshold.

- Anti-hebb guarantees neurons are decorrelated.

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

where p^2 because this is what we would expect if i and j were decorrelated. There more coactive two neurons are, the more this drives them to repulse one another.

- Standard hebbian rule

$$\Delta q_{ij} = \beta y_i(x_j - q_{ij}) \quad (3)$$

relates to sparsity fraction of neurons.

- Problems:

- Don't know how to deal with graded (i.e. non-binary) input signals. Non-discrete stuff.
- No objective function. Would like way to characterize how well system is performing.

- Led to Bruno's work: [sparse coding for graded signals](#)

- Data described by

$$I(x, y) = \sum_i a_i \phi_i(x, y) + \epsilon(x, y) \quad (4)$$

- basis decomposition of input. neuron i with activity a_i means that need feature functions ϕ to describe model. Want the [neural activities \$a_i\$ to be sparse](#).
- Constrain sparseness of a_i by imposing cost function on the activity:

$$E = \frac{1}{2} \|I - \Phi a\|^2 + \lambda \sum_i C(a_i) \quad (5)$$

where first term: preserve information and second term: I want to be sparse.

- Penalty function C shaped like really steep parabola on zero. Or could do $C = |a_i|$, v-shaped thing.

- Energy function determines dynamics of system. Want neuron activity to be expressible as a function of the input I .
- Compute coefficients a_i by gradient descent.

$$\tau \dot{a}_i = -\frac{dE}{da_i} \quad (6)$$

- Neuron i inhibited by neuron j proportionally to their functions ϕ_i inner products.
- self-inhibition of neuron back on itself makes it sparse.
- Learning rule:

$$\Delta \phi_i = -\eta \frac{\partial E}{\partial \phi_i} \quad (7)$$

$$= [I - \Phi \hat{a}] \hat{a}_i \quad (8)$$

- Abstract: A layer of simple Hebbian units connected by modifiable anti-Hebbian feed-back connections can learn to code a set of patterns in such a way that statistical dependency between the elements of the representation is reduced, while information is preserved.
- *Introduction*
 - Input-space that is our surrounding is enormous, but most inputs are highly correlated, which the brain may exploit to transform the high-dimensional pattern inputs to symbolic representations. Objects may be defined as conjunctions of highly correlated sets of components that are relatively independent of other such conjunctions¹
- *Unsupervised Learning*
 - The complexity of the mapping to be learnt \Leftarrow complexity of the input.
 - Unsupervised learning exploits statistical regularities in input to learn a more meaningful symbolic representation.
- *The Hebb Unit*
 - Simple model of cell (basically perceptron)
$$y = \begin{cases} 1 & \sum_j w_j x_j > \text{thresh} \\ 0 & \text{otherwise} \end{cases} \quad (9)$$
 - Can be thought of as pattern matching; y is maximal when weight vector = input vector pattern.
 - Hebb proposed: connection should become stronger if the two units being connected are active simultaneously: $\Delta w_j = x_j y$.
- *Competitive Learning*

¹Translation: objects are clumps of stuff that are usually found clumped together, and such that these clumps tend not to clump with other clumps.

- Out of the units receiving weighted sums of the input, only activate the unit with the largest weighted sum; suppress the output of all others.
- Results in a local, "grandmother-cell" representation.
- Limited in number of different inputs it can discriminate, and in ability to generalize.

- ***Sparse Coding***

- Distributed coding: instead, code each input state by a set of active units (rather than just one).
- Pros: combinatorics of input states increases representational capacity. Cons: situations where many units are active per input pattern, and fact that learning can be extremely slow.
- **Sparse Coding** is a compromise between distributed and local representations.

- ***Decorrelation***

- Units *within* a layer are connected by modifiable inhibitory weights, governed by an **Anti-Hebbian learning rule**: if two units in same layer are active, connection becomes more inhibitory².

²which discourages their joint activity

1.3.1 UNSUPERVISED LEARNING

• **Bruno:PCA**

- First, let's get this straight. Difference between **covariance** and **correlation**:

$$\mathbf{COV}[X, Y] \triangleq \mathbb{E}[(X - \mu_X)(Y - \mu_Y)] \quad (10)$$

$$\mathbf{CORR}[X, Y] \equiv \rho_{XY} \triangleq \frac{\text{Cov}[X, Y]}{\sigma_X \sigma_Y} \quad (\text{corr})$$

- Consider input stream \mathbf{x} that has linear pairwise correlations³ among its elements. Mathematically, correlation between elements x_i and x_j would imply that

$$\langle x_i x_j \rangle = \frac{\mathbb{E}[x_i x_j]}{\sqrt{\mathbb{E}[x_i] \mathbb{E}[x_j]}} \neq 0 \quad (11)$$

or, equivalently, that $\mathbb{E}[x_i x_j] \neq \mathbb{E}[x_i] \mathbb{E}[x_j] = 0$. Bruno is correct that linear pairwise correlations imply that $c_{ij} \neq 0$, he is *absolutely incorrect* to say that c_{ij} is an “average over many examples.” That is nothing more than academic sloppiness at its finest.

• **HKP:PCA**

- Goal: Find a set of M orthogonal vectors in data space that account for as much as possible of the data's variance. Projecting the data from original N -dimensional space onto the M -dimensional subspace spanned by these vectors then performs a **dimensionality reduction**.
- HKP actually states accurately what Bruno meant to state: The k th principal component direction is along an eigenvector direction belonging to the k th largest eigenvalue of the full **covariance matrix**

$$\langle (\xi_i - \mu_i)(\xi_j - \mu_j) \rangle \quad (12)$$

³This is exactly what is meant by eq corr, Pearson's correlation coefficient. *Linear* because “it is a measure of the linear dependence between two variables X and Y.”

- Note: I am now going to start from beginning of CH8 of HKP since I’m not understanding the stuff they are referencing FML

- **HKP Ch8: Unsupervised Hebbian Learning**

- Units need to learn patterns/correlations/categories in inputs and code the output. Units and connections display some degree of **self-organization**.
- Redundancy provides knowledge: w/o redundancy there would be no patterns to learn.

$$\text{MaxInfoPossible} - \text{InputContent} = \text{DegreeOfRedundancy} \quad (13)$$

- **PLAIN HEBBIAN LEARNING.** Context: output will be continuous-valued and DO NOT have a winner-take-all character⁴, and so the **purpose** is to measure familiarity or projecting onto principal components of input data.
- Setup: Draw at each time step an input vector ξ from (multivariate) probability distribution $P(\xi)$ that has N components⁵. Network will learn to tell us - as output - how well an input conforms to the distribution⁶
- (One linear output unit): Let V be a scalar-valued continuous output with a bunch of inputs pointing to it, with

$$V = \sum_j w_j \xi_j = \mathbf{w}^T \xi = \xi^T \mathbf{w} \quad (14)$$

- Want large (on average) $V \leftrightarrow$ more probable ξ . Why? Because then we can use the relative size of the output as a way of characterizing the sort of input it just received (see footnote 26 below). The weight update to do this is **plain Hebbian learning update**:

$$\Delta w_i = \eta V \xi_i \quad (15)$$

⁴TODO: Come back and explain why this is true, because current Brandon thought otherwise.

⁵Confusingly, here N refers to the dimension of space that each input vector lives in (usually denoted by d .)

⁶**Q:** Come back and explain why we would want a network to do this. Biological relevance/analog? **A:** You need to view it in the context of the grandmother-cell. That’s what this is all about. If a given neuron has a large linear output, then we have a good idea of what type of input went in; it was an input really similar to the weight vector. This begs the question, though: how does one determine a reasonable initialization for a given connected layer of weights to a single output? I suppose the answer is that this is the wrong question. Rather, we should interpret the outcome as resulting from a stream of particular inputs and, based on its future responses to inputs, we can determine what type of input went in. With the brain, this is like the Jennifer Aniston neuron: if that neuron fires, we can assume the person just saw something that resembled Jennifer Aniston.

where it is perhaps easier to think about the situation where $\Delta w_i = 0$ when analyzing, i.e. If $\xi_i = 0$ (which means it had nothing to do with the output), then don't increase it's weight⁷.

- Problem: \mathbf{w} grows without bound. However, suppose stable equilib exists for \mathbf{w} . This could happen for example, when considering that the update just performs $\mathbf{w} = \eta V \boldsymbol{\xi}$, where eventually $\|\mathbf{w}\| \gg \|\boldsymbol{\xi}\|$ in addition to the fact that $\boldsymbol{\xi}$ is quite likely to be along \mathbf{w} . So at equilib, expect the updates to average to 0:

$$0 = \langle \Delta w_i \rangle \quad (16)$$

$$= \left\langle \sum_j w_j \xi_j \xi_i \right\rangle \quad (17)$$

$$= \sum_j C_{ij} w_j \quad (18)$$

where the brackets are *expectation values* in the sense that

$$\langle \xi_i \xi_j \rangle = \iint_{-\infty}^{\infty} \xi_i \xi_j f_{\xi_i \xi_j}(\xi_i, \xi_j) d\xi_i d\xi_j \quad (19)$$

where f is the PDF for the two random variables in question. I suppose that, since strictly speaking \mathbf{w} isn't a random variable, that it can be pulled out along with the summation. That satisfies me for now.

- Given that $\boldsymbol{\xi}$ can be interpreted as a column vector, we have

$$C_{ij} \equiv \langle \xi_i \xi_j \rangle \quad (20)$$

$$\mathbf{C} \equiv \langle \boldsymbol{\xi} \boldsymbol{\xi}^T \rangle \quad (21)$$

Now, to be perfectly clear, this is NOT the correlation, but I am so sick and tired of caring that I'm just going to accept their absolutely incorrect definition and move on.

- Since I've read ahead, I know that the following property will be important to remember:

$$\forall \mathbf{x}, \mathbf{x}^T \mathbf{C} \mathbf{x} = \mathbf{x}^T \langle \boldsymbol{\xi} \boldsymbol{\xi}^T \rangle \mathbf{x} \quad (22)$$

$$= \langle \mathbf{x}^T \boldsymbol{\xi} \boldsymbol{\xi}^T \mathbf{x} \rangle \quad (23)$$

$$= \langle (\boldsymbol{\xi}^T \mathbf{x})^2 \rangle \quad (24)$$

- There are *only* unstable fixed points (unstable equilib) for the plain Hebbian learning procedure.
- **OJA'S RULE**. Goal: Modify plain Hebb rule such that $|\mathbf{w}| = 1$.
- Solution: Add a **weight decay** proportional to V^2 :

⁷Minor TODO: Analyze case of non-binary (i.e. continuous both pos/neg) inputs/outputs.

$$\Delta w_i = \eta V(\xi_i - V w_i) \quad (25)$$

and we see that Δw depends on the difference between the input and the back-propagated output⁸

- Informal analysis for zero-mean data: The average component of ξ along w will be zero, but since this is an algorithm depending on an unstable equilibrium, it will tend to fall along the maximal eigenvector of \mathbf{C} .
- Oja's rule chooses the direction of \mathbf{w} to maximize $\langle V^2 \rangle$.
- **Sanger's Learning Rule.** Setup: Now, instead of 1 output, have M output neurons with the hopes that they gives us the first M principal components of the input data. Architecture is ONE LAYER fully connected.
- The i th output is a linear neuron as usual given by

$$V_i = \sum_j w_{ij} \xi_j = \mathbf{w}_i^T \boldsymbol{\xi} = \boldsymbol{\xi}^T \mathbf{w}_i \quad (26)$$

- The Sanger's learning rule update for the connection *from* the j th input component *to* the i th output neuron (so we are only updating a single edge/line in the following) is

$$\Delta w_{ij} = \eta V_i \left(\xi_j - \sum_{k=1}^i V_k w_{kj} \right) \quad (27)$$

where the (converged) weight vectors to the output neurons are orthonormal and converge to the normalized eigenvectors in order of largest to smallest eigvals:

$$\mathbf{w}_i^T \mathbf{w}_j = \delta_{ij} \quad (28)$$

$$\mathbf{w}_i \rightarrow \pm \mathbf{c}^i \quad (29)$$

⁸Say 'back-propagated output' because we are subtracting what was put into the network by the resultant output *times* the connection (weight) between the input and said output. Dwelling on this *would* be overly pedantic, so move on.

- Want to learn a “dictionary” from data
- Encode input data such that it can be reconstructed from that code, where $\dim(\text{encoding}) \leq \dim(\text{input})$.
- Given N -dimensional input, build $N \times M$ dictionary⁹ (matrix) Φ where each column ϕ_i is a dictionary element with corresponding coefficient¹⁰ a_i . Want to assemble $a_i \phi_i$ into a vector of **activations**.
- **GOAL:** Minimize energy function E , defined as

$$E = \frac{1}{2} \|S - \hat{S}\|_2^2 + \lambda \sum_i^M C(a_i) \quad (30)$$

where $\hat{S} = \sum_i^M a_i \phi_i$ is for some reason called the image reconstruction. View this like a regularization procedure where the terms mean: (1) smallest difference between true image and reconstructed image (**reconstruction quality**); and (2) limit the number of **active elements**¹¹ a_i .

- Want to minimize E such that reconstructs data with fewest number of active elements, expressed as

$$\arg \min_{a, \Phi} (E) \quad (\text{argminE})$$

where I guess the double argmin means “minimize E by changing a and Φ only and then give me the values of a and Φ ”.

- Popular cost function is the ℓ_1 penalty:

$$\sum_i^M C(a_i) = \sum_i^M |a_i| \quad (31)$$

⁹M \leq N.

¹⁰Looks like $a_i \notin \Phi$

¹¹A.k.a sparsity constraints a.k.a limit activations.

- We compute coeff vector a using a “dynamic process”¹² that minimizes $\text{argmin} E$.
- Method for computing the sparse code from a given input signal S and dictionary element ϕ_i is the **Locally Competitive Algorithm**.

The model describes an activation coefficient, a_k , as the thresholded output of some model neuron’s **internal state**, u_k , which is analogous to the neuron’s membrane potential.

- Here we compute the equation for state transitions (updates) from the energy function. First, for grad descent on an individual neuron’s activity, $a_k(t)$:

$$-\frac{\partial E(t)}{\partial a_k(t)} = \sum_i^N \left[S_i \Phi_{ik} - \sum_{j \neq k}^M \Phi_{ik} \Phi_{ij} a_j \right] - a_k - \lambda \frac{\partial C(a_k)}{\partial a_k} \quad (32)$$

where the constants are S and Φ . Want system to evolve over time to produce optimal set of activations $a(t)$.

- Meaning of ϕ_k . Associated with k th (output?) neuron. Indicates the connection strength [between that neuron and] each pixel in the input.

In this model, we are going to find a sparse code for one patch of an image at a time, so that all M neurons are connected to the same image patch, S .

¹²Okay well what the fuck is it?

Nearby (similar) outputs corresponding to nearby (similar) input patterns. Such a map (similar inputs \rightarrow similar outputs) is a **feature map**. The conventional case: 2 continuous-valued inputs x and y map (fully-connected) to a two-dimensional x,y grid. Want nearby input values (in the actual euclidean sense) (x, y) to be mapped closely in the output 2D grid.

Kohonen's Algorithm implements the self-organizing (feature) map by using competitive learning, where now we update weights going to the *neighbors* of the winning unit as well as those of the winning unit itself.

- Setup: N continuous-valued inputs ξ_1 to ξ_N , defining a point ξ in N -dimensional space. Outputs O_i are arranged in (typically) a 1-D or 2-D array fully connected via w_{ij} to the inputs.
- A competitive learning rule is used, choosing output O_i^* as winner, determined by

$$|w_i^* - \xi| \leq |w_i - \xi| \quad (\text{for all } i) \quad (33)$$

- The **Kohonen Learning Rule** is

$$\Delta w_{ij} = \eta \Lambda(i, i^*) (\xi_j - w_{ij}) \quad (34)$$

where $\Lambda(i, i^*)$ is the **neighborhood function**, equal to 1 for $i = i^*$ and falls off with distance $|\mathbf{r} - \mathbf{r}_i^*|$.

- A typical choice for $\Lambda(i, i^*)$ is

$$\Lambda(i, i^*) = \exp\left(-\frac{|\mathbf{r} - \mathbf{r}_i^*|^2}{2\sigma^2}\right) \quad (35)$$

where σ is width parameter that *is gradually decreased*. Apparently $\eta(t) \propto t^{-\alpha}$ where $0 < \alpha \leq 1$ is a good choice.

Locally Linear Embedding

Table of Contents Local

Written by Brandon McKinzie

LLE is an unsupervised learning algorithm for dimensionality reduction. Similar to PCA and MDS¹³, LLE is called an *eigenvector method*. The basic idea is illustrated below in figure 1.

The **LLE algorithm**:

1. Compute the neighbors of each data point, X_i .
2. Compute the weights W_{ij} that best reconstruct each X_i from its neighbors, minimizing the cost in

$$ReconErr(W) = \sum_i |X_i - \sum_j W_{ij} X_j|^2 \quad (36)$$

by constrained linear fits.

3. Compute the Y_i reconstructed by the weights W_{ij} , minimizing the quadratic form in

$$\Phi(Y) = \sum_i |Y_i - \sum_j W_{ij} Y_j|^2 \quad (37)$$

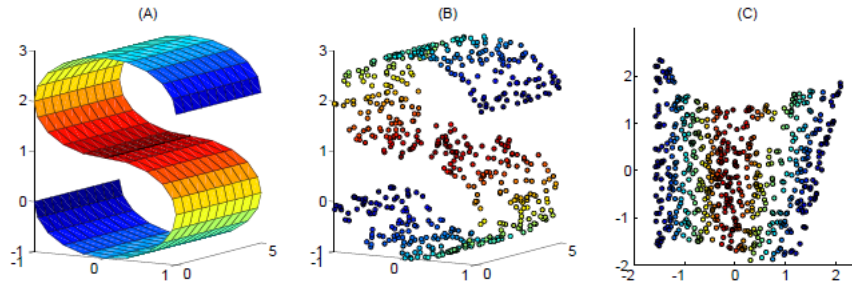


Figure 1: (A) Multidimensional sampling distribution with clear underlying manifold representation. (B) Points that were sampled. (C) The neighborhood-preserving mapping discovered by LLE.

¹³Multidimensional scaling

by its bottom nonzero eigenvectors.

Some intuition/overview of the algorithm. We expect each X_i and its neighbors to lie on or close to a **locally linear** patch of the manifold. We characterize these patches by linear coefficients W_{ij} that reconstruct each X_i from its neighbors. As seen in eq. 36, the reconstructed point X_i is given by $\sum_j W_{ij} X_j$.

Computing/analyzing the weights W_{ij} . Minimize eq 36 subject to

$$\forall X_j \notin \text{Neighbors}(X_i) : W_{ij} = 0 \quad (38)$$

$$\sum_j W_{ij} = 1 \quad (39)$$

where the optimal weights are found by solving a least squares problem. Note that for a given data point, *the weights are invariant to rotations, rescalings, and translations of that data point and its neighbors.*¹⁴ If the data lie on some nonlinear manifold of $d \ll D$, then there exists a *linear mapping* (approx) from the high-D coordinates of each neighborhood to global ('internal') coordinates on the manifold. Lucky for us, W can also do this!¹⁵

Explanation of eqs. 36 37. Note that eq. 36 is minimized over the W_{ij} , while equation 37 is minimized over the Y_i . In English: We first want the weights W that reconstruct each X_i by its neighbors in the high-D space. Then, we want the low- d coordinates Y_i , representing the global coordinates on the manifold, that correspond to each X_i from the original space. **How it is minimized:**

it can be minimized by solving a sparse $N \times N$ eigenvector problem, whose bottom d non-zero eigenvectors provide an ordered set of orthogonal coordinates centered on the origin.

¹⁴In other words, since the weights just characterize the local patch of the given data point, that patch shouldn't change if we shift the data, rotate it, or scale it. The neighboring points should remain the same.

¹⁵In particular, the same weights W_{ij} that reconstruct the i th data point in D dimensions should also reconstruct its embedded manifold coordinates in d dimensions

Implementation of algorithm. Only one free parameter: number of neighbors per data point K . W_{ij} and Y_i are computed by 'standard linear algebra'.

Recurrent Neural Networks

Table of Contents Local

Written by Brandon McKinzie

Lab 6 Overview. Briefly goes over how we can corrupt some number of bits and reconstruct a desired image [with hopfield nets]. Unfortunately, can get “spurious basins of attractions.” Pushing down on some region of landscape causes pushing up of some other region. Want to carve energy landscape so that we push down only where we want.

Bump circuits and ring attractors. Want family of solutions (e.g. a line) that solutions drawn to (called line attractors). Head-direction neurons¹⁶ look like an internal compass for animals; encode direction of head in *world coordinate system*. Different dots represent a single neuron’s firing rate at different relative head directions. **Ring attractors:** population of neurons that with bumps that are stable (?). Convergence/stability because T_{ij} matrix is symmetric. Symmetric = fixed stable; Asymmetric =

Bruno shows simulation:

- 32 neurons where bar is activity of neuron.
- Start with random symmetric weight hopfield net.
- Eventually weights converge to gaussian-like bump; an equipotential pattern.
- If we add small asymmetry (gamma) to weights, then population (bump) would shift. Bump change is shifting position, and when the asymmetry stops (we stop moving our head) the population stays fixed. In English: moving head causes bump to move but when we stop moving, they stay put.
- **For more:** Read “catcher and zong” paper. I misspelled that.

¹⁶Literally referring to direction of [e.g. some animal’s] head

[Enter guest lecturer Alex Anderson] **Recurrent Neural Networks:**

→ Starts with handwriting network.

→ RNNs good for sequence prediction tasks with “long-term dependencies.”

Backprop Review. Blobs do activation computation and transformers do propagations. Note: A^t is target output values.

Problem to Solve. Feed net a bunch of sentences and have it fill in the blank somewhere, based on the previous info it was fed. Mad libs. Have network understand particular frame of movie by exploiting context; just showing it a bunch of frames isn’t enough/good approach.

RNN loops/Notation. Feed *time sequence* x_t to block A . Two figures in this slide are different reps of same thing; instructor prefers the right fig. H_k is hidden state we want to predict¹⁷. f can be some nonlinearity like *tanh*. In RNNs, cost function typically broken up over time; so C_k is cost at timestep k . Usually want hidden state to *summarize* the past. Hidden state traces out a trajectory over time [wut].

Unroll a RNN. Can basically turn RNN into a linearized hidden markov chain, where time proceeds to the right. Total cost is given by cost at each time step.

Long-term Dependencies. Shows toy model. Imagine ur an ant walking along graph. Given string of nodes, predict next letter each timestep [solve the question mark in slide]. Don’t necessarily want/need whole past as input. Want to remember past [hidden] states, but they usually get overwritten; want to save it more efficiently. Key: want to make function simple, give the network parameterization.

Exploding/vanishing gradients. Local dependencies easy to learn.

→ Once we get to B, want network to output a U.

→ To learn, errors need to propagate back [in time], so we can change the weights that started the error: gradient of cost at timestep k with respect to initial weights using chain rule. Basically a product of k matrices.

→ If k large and matrices have eigvals less than 1, gradients *vanish*. If eigvals above 1, gradients *explode*. So what we want is for eigvals to be very near 1.

→ **Todo:** lookup relationship between eigval magnitudes and determinant.

¹⁷Analogy to hopfield: H is like hopfield B. X is like external I in hopfield.

Solution: Multiplicative Gating. Helps protect hidden state. MultGate can be either 0 or 1, and we multiply the hidden state by that value; if we 0 lose the hidden state; if 1 we keep the hidden state. Since binary functions not smooth/differentiable, continuous gating is better. [slide note: top row is w/o multgate, lower row is with multgate]. Key equation:

$$c_t = f_t \odot c_{t-1} + i_t \odot j_t$$

where \odot is elementwise product.

Note: This is in TensorFlow now.

Hopfield Networks Handout

Table of Contents Local

Written by Brandon McKinzie

Energy Function. The following governs the dynamic of pairwise recurrently connected networks.

$$E = \frac{1}{2} \sum_i \sum_{j \neq i} T_{ij} V_i V_j \quad (40)$$

For symmetric weights $T_{ij} = T_{ji}$, consider the change in energy ΔE resulting from making a positive change to V_k ¹⁸

$$\Delta E = -\Delta V_k \sum_{i \neq k} T_{ki} V_i \quad (41)$$

which will be *negative* if both ΔV_k and the sum are positive, thus decreasing the overall energy (good). Conversely, if sum is negative, we should decrease value of V_k . **Critical assumption:** Symmetric $T_{ij} = T_{ji}$. Without this assumption, impossible to show the system will have fixed points.

For a network with symmetric connections though, the dynamics will converge to so-called **basins of attraction**.

Setting the Weights. Goal: store pattern \mathbf{V}^α as basin of attraction in network. One approach: the **Hebbian prescription** $T_{ij} = V_i^\alpha V_j^\alpha$.

→ **Single memory storage.** Now, the summed input sent to, say, the i th unit in response to some \mathbf{V}^β will be given by

$$U_i = V_i^\alpha \sum_{j \neq i} V_j^\alpha V_j^\beta \quad (42)$$

and thus if $\mathbf{V}^\alpha = \mathbf{V}^\beta$, U_i won't flip sign and the networks stays put.

→ **Multiple memories.** Now, need to form as many basins of attractions as memories we want stored. Set weights with a superposition over each desired *memory* \mathbf{V}^α : $T_{ij} = \sum_\alpha V_i^\alpha V_j^\alpha$, and the corresponding response of the i th neuron is

$$U_i = \sum_\alpha V_i^\alpha \sum_{j \neq i} V_j^\alpha V_j^\beta \quad (43)$$

¹⁸If it is -1, change to +1, else just keep where it is.

Capacity for a Hopfield Network. If the patterns to store (memories) have *few elements in common*, then cross terms $\sum_{j \neq i} V_j^\alpha V_j^\beta$ tend to zero for $\alpha \neq \beta$ (since each V_j^α is ± 1 and a random average over ± 1 is zero) and U_i won't change. As we store more patterns which are *similar*, memories degrade and basins gone from desired locations. This **capacity** for Hopfield is $\approx 15\%$ of the number of neurons in network¹⁹.

¹⁹Assuming the stored patterns are relatively dissimilar.

Mixture of Gaussians and EM Algorithm

Table of Contents Local

Written by Brandon McKinzie

(MOG) Model Assumptions. Assume each $x^{(i)}$ generated by sampling $z^{(i)} \sim \text{Multinom}(\phi)$ ²⁰ and then positing that $x^{(i)}$ was drawn from the Gaussian associated with $z^{(i)}$ (so k possible Gaussians since z could be one of k classes), i.e.

$$x^{(i)} | z^{(i)} = j \sim \mathcal{N}(\mu_j, \Sigma_j) \quad (44)$$

Goal. Estimate the parameters of our model, ϕ, μ, Σ . We can do this by writing the likelihood of our data (m data points):

$$\ell(\phi, \mu, \Sigma) = \sum_{i=1}^m \log \left[\sum_{z^{(i)}=1}^k p(x^{(i)} | z^{(i)}; \mu, \Sigma) p(z^{(i)}; \phi) \right] \quad (45)$$

where the inner sum comes from using Bayes rule on $p(x^{(i)})$. Note that *we don't know the $z^{(i)}$ for each data point*, that information is hidden. This means we can't get closed-form solutions by doing MLE/taking derivatives as usual.

EM Algorithm. Purpose: allow us to move forward given we can't do the standard MLE approach (since we don't know the values of each $z^{(i)}$.)

1. **E-step.** Estimate the values of the $z^{(i)}$ s by evaluating

$$\text{For each } i, j, \text{ set } w_j^{(i)} := p(z^{(i)} = j | x^{(i)}; \phi, \mu, \Sigma) \quad (46)$$

with Bayes' rule, using whatever the current values of the estimated parameters are.

²⁰where k -vector ϕ elements $\phi_j = P(z^{(i)} = j)$

2. **M-step.** Update parameters via standard MLE approach.

$$\phi_j := \frac{1}{m} \sum_{i=1}^m w_j^{(i)} \quad (47)$$

$$\mu_j := \frac{\sum_{i=1}^m w_j^{(i)} x^{(i)}}{\sum_{i=1}^m w_j^{(i)}} \quad (48)$$

$$\Sigma_j := \frac{\sum_{i=1}^m w_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^m w_j^{(i)}} \quad (49)$$

Some properties of the EM-algorithm. It is similar to the K-means algorithm, and like K-means, it is also prone to local minima, so reinitializing at several different initial parameters may be a good idea.

Boltzmann Machines

Table of Contents Local

Written by Brandon McKinzie

1.10.1 LECTURE SLIDES

[Here, I use my own notation that actually remains consistent/intuitive...]

The energy, taken as a sort of average over all pairwise connections and the corresponding network-wide probability is

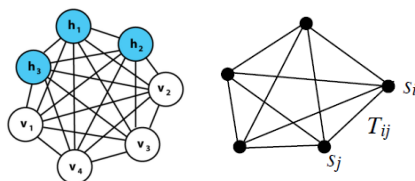
$$E(\mathbf{s}) = -\frac{1}{2} \sum_{ij} w_{ij} s_i s_j \quad (50)$$

$$P(\mathbf{s}) = \frac{1}{Z} e^{-\beta E(\mathbf{s})} \quad (51)$$

1.10.2 HKP CHAPTER 7.1

Structure. The units S_i are divided into

1. **Visible units.** These may be further divided into, e.g., input and output units.
2. **Hidden units.** No connection to the outside world.



The units are **stochastic**, meaning

$$S_i = \begin{cases} +1 & \text{with probability } g(h_i) \\ -1 & \text{with probability } 1 - g(h_i) \end{cases} \quad (52)$$

$$h_i = \sum_j w_{ij} S_j \quad (53)$$

$$g(h) = \frac{1}{1 + \exp(-2\beta h)} \quad (54)$$

where $\beta = 1/T$ and T is the “temperature.” **Goal:** adjust the weights w_{ij} to give the states of the *visible* units a particular desired probability distribution. This differs from a Hopfield network in that now we have hidden units. With hidden units, we can specify *higher-order correlations* between units²¹.

1.10.3 BOLTZMANN MACHINES - AIFH

A Boltzmann machine is essentially a fully connected, two-layer neural network; one visual layer and one hidden layer. **Restricted** Boltzmann machines are not fully connected; all hidden neurons are connected to each visible neuron and vice versa, but there are no connections between neurons of the same layer.

- Boltzmann machines are a generative model.
- The values presented to the visible neurons of a Boltzmann machines, when considered with the weights, specify a probability that the hidden neurons will assume a value of 1, as opposed to 0.

Differences with Hopfield networks.

- Hopfield networks suffer from recognizing false patterns.
- BM can store a greater capacity of patterns than HN.
- HN require the input patterns to be uncorrelated.
- BM can be stacked to form layers.

²¹Whereas, in hopfield networks, we can do no more than specify all the $\langle S_i \rangle$ and $\langle S_i S_j \rangle$

Independent Component Analysis

Table of Contents Local

Written by Brandon McKinzie

Sparse Coding Review. The goal is to represent input x as vector of sparse coefficients s , where their relationship is given in the form

$$\mathbf{x} = \mathbf{A}\mathbf{s} + \mathbf{n} \quad (55)$$

- \mathbf{x} : Data matrix with shape n by d .
- \mathbf{A} : Feature matrix (contains the basis functions) of shape n by $k > n$. Here, k is the number of sparse coefficients (shape of \mathbf{s}).
- \mathbf{s} : The sparse coefficient matrix. Recall that the equation above is sometimes written as a sum over the basis functions ϕ .

$$\mathbf{x} = \left(\sum_k \phi_k(\mathbf{x}) s_k \right) + \mathbf{n} \quad (56)$$

We say “sparse” because $s_i = 0$ a lot, and $\text{shape}(\mathbf{s}) > \text{shape}(\mathbf{x})$.

→ \mathbf{n} : Gaussian noise. Entries are typically much smaller than entries of $\mathbf{A}\mathbf{s}$.

- **Model distribution.**

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{s})p_s(\mathbf{s})d\mathbf{s} \quad \text{where} \quad (57)$$

$$p(\mathbf{x}|\mathbf{s}) \propto e^{-\frac{|\mathbf{x}-\mathbf{A}\mathbf{s}|^2}{2\sigma_n^2}} \quad \text{and} \quad (58)$$

$$p_s(\mathbf{s}) \propto e^{-\sum_i C(s_i)} \quad \leftrightarrow \quad \mathbf{C} = -\log(p(\mathbf{s})) \quad (59)$$

- **Learning Rule.**

$$\Delta \mathbf{A} \propto \frac{\partial}{\partial \mathbf{A}} \langle \log p(\mathbf{x}) \rangle \quad (60)$$

$$\Delta \mathbf{A} \propto \left\langle \int [\mathbf{x} - \mathbf{A}\mathbf{s}] \mathbf{s}^T p(\mathbf{s}|\mathbf{x}) d\mathbf{s} \right\rangle \quad (\text{ANALYTIC}) \quad (61)$$

$$\Delta \mathbf{A} \propto \langle [\mathbf{x} - \mathbf{A}\hat{\mathbf{s}}] \hat{\mathbf{s}}^T \rangle \quad (\text{IN PRACTICE}) \quad (62)$$

where, $\hat{\mathbf{s}}$ represents a single sample at the posterior maximum:

$$\hat{\mathbf{s}} = \arg \max_{\mathbf{s}} p(\mathbf{s}|\mathbf{x}) \quad (63)$$

$$= \arg \min_{\mathbf{s}} [-\log(p(\mathbf{s}|\mathbf{x}))] \quad (64)$$

$$= \arg \min_{\mathbf{s}} \left[\frac{\lambda_n}{2} |\mathbf{x} - \mathbf{A}\mathbf{s}|^2 + \sum_i C(s_i) \right] \quad (65)$$

$$\nabla_{\mathbf{s}} \hat{\mathbf{s}} \propto \lambda_n \mathbf{A}^T [\mathbf{x} - \mathbf{A}\mathbf{s}] - C'(\mathbf{s}) \quad (66)$$

$$\triangleq \lambda_n [\mathbf{b} - \mathbf{G}\mathbf{s}] - \mathbf{z}(\mathbf{s}) \quad (67)$$

Independent Component Analysis. Special case where (1) \mathbf{A} is square and full rank, and (2) the noise $\mathbf{n} = 0$. Now we have model

$$\mathbf{x} = \mathbf{A}\mathbf{s} \quad \rightarrow \quad \mathbf{s} = \mathbf{A}^{-1}\mathbf{x} \quad (68)$$

Learning rule.

$$\Delta \mathbf{A} \propto \langle [\mathbf{x} - \mathbf{A}\hat{\mathbf{s}}] \hat{\mathbf{s}}^T \rangle \quad (69)$$

$$\Delta \mathbf{A} \propto \mathbf{A} \langle \mathbf{z}(\mathbf{s}) \mathbf{s}^T \rangle - \mathbf{A} \quad (70)$$

Equations for Different Priors

$$\text{[LAPLACE]} \quad P(s_i) \propto e^{-|s_i|} \quad \leftrightarrow \quad z_i = \text{sign}(s_i) \quad (71)$$

$$\text{[CAUCHY]} \quad P(s_i) \propto \frac{1}{1 + s_i^2} \quad \leftrightarrow \quad z_i = \frac{2|s_i|}{1 + s_i^2} \quad (72)$$

$$\text{[GAUSS]} \quad P(s_i) \propto e^{-s_i^2} \quad \leftrightarrow \quad z_i = |s_i| \quad (73)$$

Algorithm summary/procedure.

1. Initialize square matrix A .
2. Until A converges, do:
 - Compute source vector via $\mathbf{s} = \mathbf{A}^{-1}\mathbf{x}$.
 - Compute

$$\mathbf{z} = \nabla_{\mathbf{s}} C(\mathbf{s}) \quad (74)$$

$$= \nabla_{\mathbf{s}} [-\log(p(\mathbf{s}))] \quad (75)$$

$$= -\sum_i \nabla_{\mathbf{s}} \log(p_s(s_i)) \quad (76)$$

- Update

$$\Delta A \propto A \langle \mathbf{z}(\mathbf{s}) \mathbf{s}^T \rangle - A \quad (77)$$

1.11.1 ICA - ANDREW NG - CS 229

Cocktail Party. There are n people talking at a cocktail party, and we've placed n microphones in the room to see if we can separate out the original n speakers' speech signals. We observe

$$\mathbf{x} = \mathbf{A}\mathbf{s} \quad (78)$$

- $\mathbf{x}^{(i)}$ denotes the n -dimensional vector of our microphone recordings at time i .
- $\mathbf{s}^{(i)}$ denotes the n -dimensional vector of each speakers' output at time i .
- \mathbf{A} be the unknown square **mixing matrix**. **Goal:** Find the *unmixing matrix* $\mathbf{W} = \mathbf{A}^{-1}$. Then we can recover the generated sources via $\mathbf{s}^{(i)} = \mathbf{W}\mathbf{x}^{(i)}$.

$$\mathbf{W} = \begin{bmatrix} - & w_1^T & - \\ - & w_2^T & - \\ & \vdots & \\ - & w_n^T & - \end{bmatrix} \quad \longrightarrow \quad s_j^{(i)} = w_j^T \mathbf{x}^{(i)} \quad (79)$$

[ELI5] “People say stuff $\mathbf{s}^{(i)}$, but it gets all mixed up so we hear stuff $\mathbf{x}^{(i)}$. We want to know how they related. Luckily, we can just unmix the mixed stuff.”

Ambiguities. The (1) order and (2) scaling of the n “ w_i ” vectors in W is ambiguous. The sign of s_i is irrelevant (sounds the same on a speaker). **Key point:** The aforementioned ambiguities are the *only* ambiguities, so long as the sources s_i are *non-Gaussian*. Basically, this is due to things like rotational invariance that can be present in Gaussians.

ICA Algorithm.

1. Let each source s_i have density $p_s(s_i)$. Then the joint distribution $p(s)$ is the probability of hearing all independent sources s_i . From previous results, we know that this implies a density for $x = As = W^{-1}s$ (below).

$$p(s) = \prod_{i=1}^n p_s(s_i) \quad \implies \quad p(x) = \prod_{i=1}^n p_s(w_i^T x) \cdot |W| \quad (80)$$

FINAL PROJECT

CONTENTS

2.1	WYGIWYS - A Visual Markup Decompiler	32
-----	--	----

WYGIWYS - A Visual Markup Decompiler

Table of Contents Local

*Written by Brandon McKinzie***Inputs/Outputs.**

- **Input:** $\mathbf{x} \in \mathcal{X}$, consists of an image e.g. $\mathbb{R}^{H \times W}$ for grayscale images.
- **Output:** $\mathbf{y} \in \mathcal{Y}$, where $\mathbf{y} = \langle y_1, \dots, y_C \rangle$ contains C tokens in the markup language.
- **Example:** Below, should we feed the input as the image of the equation on the left, the output would be the vector on the right.

$$\rho = \sum_{\alpha > 0} \alpha \quad \langle rho, =, sum, alpha, >, 0, alpha \rangle$$

At training time, we assume we are given a sequence of input images \mathbf{x} and ground-truth \LaTeX labels \mathbf{y}

$$\left((\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots, (\mathbf{x}^{(J)}, \mathbf{y}^{(J)}) \right)$$

and at test time, given raw input image \mathbf{x} , we predict its corresponding \LaTeX source code, and output the \hat{x} after compiling our prediction, then comparing \hat{x} with x .

The Model.