

The Hebbian Descent Learning Rule

Laurenz Wiskott



Jan Melchior



INSTITUT
FÜR
NEUROINFORMATIK

RUHR
UNIVERSITÄT
BOCHUM

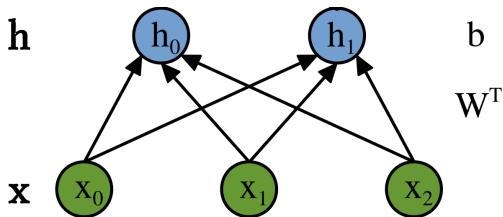
RUB

supported by:
Ruhr-University Bochum

Single layer network

Single layer network: $\mathbf{h} = \phi(\mathbf{a}) = \phi(\mathbf{W}^T(\mathbf{x} - \boldsymbol{\mu}) + \mathbf{b})$

Input \mathbf{x} , offset $\boldsymbol{\mu}$, bias \mathbf{b} , weights \mathbf{W} , activation function ϕ , output \mathbf{h} .



Gradient descent

Single layer network: $\mathbf{h} = \phi(\mathbf{a}) = \phi(\mathbf{W}^T(\mathbf{x} - \boldsymbol{\mu}) + \mathbf{b})$

Input \mathbf{x} , offset $\boldsymbol{\mu}$, bias \mathbf{b} , weights \mathbf{W} , activation function ϕ , output \mathbf{h} .

Loss function: $\mathcal{L}(\mathbf{t}, \mathbf{h})$ (with target values \mathbf{t})

Gradient descent:

$$\begin{aligned}\Delta_{GD} \mathbf{W} &= -\eta(\mathbf{x} - \boldsymbol{\mu}) \left(\frac{\partial \mathcal{L}(\mathbf{t}, \mathbf{h})}{\partial \mathbf{h}} \odot \phi'(\mathbf{a}) \right)^T \\ \Delta_{GD} \mathbf{b} &= -\eta \left(\frac{\partial \mathcal{L}(\mathbf{t}, \mathbf{h})}{\partial \mathbf{h}} \odot \phi'(\mathbf{a}) \right)\end{aligned}$$

\odot indicates element-wise multiplication.

Hebbian descent

Single layer network: $\mathbf{h} = \phi(\mathbf{a}) = \phi(\mathbf{W}^T(\mathbf{x} - \boldsymbol{\mu}) + \mathbf{b})$

Input \mathbf{x} , offset $\boldsymbol{\mu}$, bias \mathbf{b} , weights \mathbf{W} , activation function ϕ , output \mathbf{h} .

Loss function: $\mathcal{L}(\mathbf{t}, \mathbf{h})$ (with target values \mathbf{t})

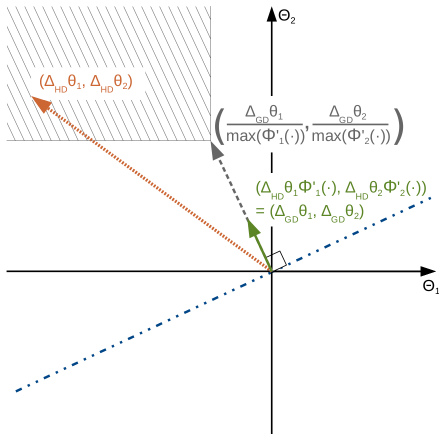
Hebbian descent:

$$\begin{aligned}\Delta_{GD} \mathbf{W} &= -\eta(\mathbf{x} - \boldsymbol{\mu}) \left(\frac{\partial \mathcal{L}(\mathbf{t}, \mathbf{h})}{\partial \mathbf{h}} \odot \phi'(\mathbf{a}) \right)^T \\ \Delta_{GD} \mathbf{b} &= -\eta \left(\frac{\partial \mathcal{L}(\mathbf{t}, \mathbf{h})}{\partial \mathbf{h}} \odot \phi'(\mathbf{a}) \right)\end{aligned}$$

\odot indicates element-wise multiplication.

- ▶ Hebbian descent differs from gradient descent by dropping the derivative $\phi'(\mathbf{a})$.
- ▶ For non-negative $\phi'(\mathbf{a})$, Hebbian descent points within 90° of the gradient.
- ▶ Hebbian descent keeps going even if $\phi'(\mathbf{a}) = 0$.

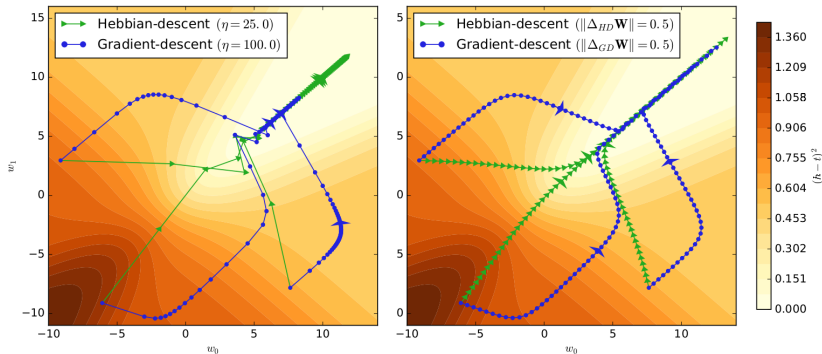
Hebbian vs gradient descent



θ stands for W or b .

- ▶ Hebbian descent is confined to the quadrant at $\Delta_{GD} \theta_i / \max(\Phi'_i(\cdot))$.

Hebbian vs gradient descent



Two input units, one output unit, trained on $\{(1, 0), (0, 1), (1, 1)\} \rightarrow \{(0), (0), (1)\}$ with L2 loss for 50 full-batch updates (wider arrowheads at 25 updates) with hand tuned learning rates avoiding oscillations. Natural step size (left), constant step size (right). Equivalent results with online learning.

- ▶ In this example, Hebbian descent finds a more direct path and converges faster.

Hebbian descent loss

Gradient descent:
$$\Delta_{GD} \mathbf{W} = -\eta(\mathbf{x} - \boldsymbol{\mu}) \left(\frac{\partial \mathcal{L}(\mathbf{t}, \mathbf{h})}{\partial \mathbf{h}} \odot \phi'(\mathbf{a}) \right)^T$$

Hebbian descent as gradient descent (for Φ with $\Phi' > 0$)

$$\Delta_{HD} \mathbf{W} = -\eta(\mathbf{x} - \boldsymbol{\mu}) \frac{\partial \mathcal{L}(\mathbf{t}, \mathbf{h})}{\partial \mathbf{h}}^T = -\eta(\mathbf{x} - \boldsymbol{\mu}) \left(\frac{\partial \mathcal{H}(\mathbf{t}, \mathbf{h})}{\partial \mathbf{h}} \odot \phi'(\mathbf{a}) \right)^T$$

on different loss
$$\mathcal{H}(\mathbf{t}, \mathbf{h}) = \sum_j^M \int \frac{\partial \mathcal{L}(\mathbf{t}, \mathbf{h})}{\partial \mathbf{h}} \odot \phi'(\mathbf{a}) dh_j$$

\odot indicates element-wise multiplication; \oslash indicates element-wise division.

- ▶ For instance, L2 loss with sigmoid leads to cross entropy loss:

$$\mathcal{L}(t_j, h_j) = \frac{1}{2}(h_j - t_j)^2 \quad \Rightarrow \quad \mathcal{H}(t_j, h_j) = -t_j \ln(h_j) - (1 - t_j) \ln(1 - h_j)$$

- ▶ Hebbian descent generalizes this to arbitrary losses and activation functions.

Hebbian descent losses

$h_j = \phi(a_j)$	$\frac{\mathcal{E}(t_j, h_j)}{\phi'(a_j)} = \frac{h_j - t_j}{\phi'(a_j)}$	$\mathcal{L}_{HD}(t_j, h_j) = \int \frac{\mathcal{E}(t_j, h_j)}{\phi'(a_j)} dh_j$
Linear a_j	$h_j - t_j$	Squared error loss $\frac{1}{2} (h_j - t_j)^2$
Sigmoid $\frac{1}{1 + \exp(-a_j)}$	$\frac{h_j - t_j}{h_j(1-h_j)}$	Cross entropy loss $-t_j \ln(h_j) - (1 - t_j) \ln(1 - h_j)$
Softmax (Bridle, 1990) $\frac{\exp(a_j)}{\sum_k \exp(a_k)}$	$\frac{h_j - t_j}{h_j(\delta_{j,k} - h_{k,j})}$	Cross entropy loss $-t_j \ln(h_j) - (1 - t_j) \ln(1 - h_j)$
Scaled Hyperbolic Tangent $\alpha \tanh(a_j)$	$\frac{h_j - t_j}{\alpha - \frac{1}{\alpha} h_j^2}$	$-\frac{1}{2}(\alpha + t_j) \ln(\alpha + h_j) - \frac{1}{2}(\alpha - t_j) \ln(\alpha - h_j)$
Approx. Step ($\alpha \rightarrow 0$) $\begin{cases} \alpha a_j \\ \alpha a_j + \beta \end{cases}$	$\begin{cases} h_j - t_j \\ \alpha \end{cases}$	$\begin{cases} \frac{1}{2\alpha} (h_j - t_j)^2 & \text{for } a_j < 0 \\ \text{for } a_j \geq 0 \end{cases}$
Leaky Rectifier. (Maas et al., 2013) $\begin{cases} \alpha a_j \\ a_j \end{cases}$	$\begin{cases} \frac{1}{\alpha} (h_j - t_j) \\ h_j - t_j \end{cases}$	$\begin{cases} \frac{1}{2\alpha} (h_j - t_j)^2 & \text{for } a_j < 0 \\ \frac{1}{2} (h_j - t_j)^2 & \text{for } a_j \geq 0 \end{cases}$
Scaled Exp. Linear (Klambauer et al., 2017) $\begin{cases} \lambda \alpha (\exp(a_j) - 1) \\ \lambda a_j \end{cases}$	$\begin{cases} \frac{h_j - t_j}{h_j - \lambda \alpha} \\ \frac{1}{\lambda} (h_j - t_j) \end{cases}$	$\begin{cases} h_j - (t_j + \lambda \alpha) \log(h_j + \lambda \alpha) & \text{for } a_j < 0 \\ \frac{1}{2\lambda} (h_j - t_j)^2 & \text{for } a_j \geq 0 \end{cases}$

$a_j = t_j$	$\frac{h_j - t_j}{h_j(1-h_j)}$	$C_{\text{CE}}(t_j, h_j) = -\int \frac{h_j - t_j}{h_j(1-h_j)} dh_j$
Linear	$h_j - t_j$	Squared error loss $\frac{1}{2} (h_j - t_j)^2$
Sigmoid	$\frac{h_j - t_j}{h_j(1-h_j)}$	Cross entropy loss $-t_j \ln(h_j) - (1 - t_j) \ln(1 - h_j)$
Softmax	$\frac{h_j - t_j}{h_j(\delta_{j,k} - h_{k,j})}$	Cross entropy loss $-t_j \ln(h_j) - (1 - t_j) \ln(1 - h_j)$
Approx. Step ($\alpha \rightarrow 0$)	$\begin{cases} h_j - t_j \\ \alpha \end{cases}$	$\begin{cases} \frac{1}{2\alpha} (h_j - t_j)^2 & \text{for } a_j < 0 \\ \text{for } a_j \geq 0 \end{cases}$
Leaky Rectifier	$\begin{cases} \frac{1}{\alpha} (h_j - t_j) \\ h_j - t_j \end{cases}$	$\begin{cases} \frac{1}{2\alpha} (h_j - t_j)^2 & \text{for } a_j < 0 \\ \frac{1}{2} (h_j - t_j)^2 & \text{for } a_j \geq 0 \end{cases}$
Scaled Exp. Linear	$\begin{cases} \frac{h_j - t_j}{h_j - \lambda \alpha} \\ \frac{1}{\lambda} (h_j - t_j) \end{cases}$	$\begin{cases} h_j - (t_j + \lambda \alpha) \log(h_j + \lambda \alpha) & \text{for } a_j < 0 \\ \frac{1}{2\lambda} (h_j - t_j)^2 & \text{for } a_j \geq 0 \end{cases}$

Hebbian descent with L2 loss

Single layer network: $\mathbf{h} = \phi(\mathbf{a}) = \phi(\mathbf{W}^T(\mathbf{x} - \boldsymbol{\mu}) + \mathbf{b})$

Input \mathbf{x} , offset $\boldsymbol{\mu}$, bias \mathbf{b} , weights \mathbf{W} , activation function ϕ , output \mathbf{h} .

Loss function: $\mathcal{L}(\mathbf{t}, \mathbf{h}) = \frac{1}{2} \|\mathbf{h} - \mathbf{t}\|^2; \quad \frac{\partial \mathcal{L}(\mathbf{t}, \mathbf{h})}{\partial \mathbf{h}} = (\mathbf{h} - \mathbf{t})$

Hebbian descent:

$$\Delta_{HD} \mathbf{W} = \underbrace{\eta(\mathbf{x} - \boldsymbol{\mu})\mathbf{t}^T}_{\text{sup. Hebb}} - \underbrace{\eta(\mathbf{x} - \boldsymbol{\mu})\mathbf{h}^T}_{\text{unsup. Hebb}}$$
$$\Delta_{HD} \mathbf{b} = \eta\mathbf{t} - \eta\mathbf{h}$$

- ▶ Difference between supervised and unsupervised Hebbian learning makes Hebbian descent stable, no learning if the desired output is achieved.
- ▶ Thus, no normalization is required.
- ▶ Thus, gradient descent benefits from multiple presentations of patterns (in contrast to Hebbian learning and the covariance rule).

Hebbian descent includes Hebb's and covariance rule

Hebbian descent with L2 loss:

$$\Delta_{HD,sup.} \mathbf{W} = \underbrace{\eta(\mathbf{x} - \boldsymbol{\mu})\mathbf{t}^T}_{\text{sup. Hebb}} - \underbrace{\eta(\mathbf{x} - \boldsymbol{\mu})\mathbf{h}^T}_{\text{unsup. Hebb}} \quad (1)$$

For $\boldsymbol{\mu} = 0$ we get Hebb's rule:

$$\begin{aligned} & \Delta_{HD,sup.} \mathbf{W} \\ &= \eta \mathbf{x} \mathbf{t}^T \\ &= \Delta_{Hebb} \mathbf{W} \end{aligned}$$

For $\boldsymbol{\mu} = \langle \mathbf{x} \rangle$ we get the covariance rule:

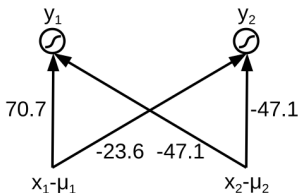
$$\begin{aligned} & \langle \Delta_{HD,sup.} \mathbf{W} \rangle \\ &= \eta \langle (\mathbf{x} - \langle \mathbf{x} \rangle) \mathbf{t}^T \rangle \\ &= \eta \langle \mathbf{x} \mathbf{t}^T \rangle - \langle \mathbf{x} \rangle \langle \mathbf{t} \rangle^T - \underbrace{\langle \mathbf{x} \rangle \langle \mathbf{t} \rangle^T + \langle \mathbf{x} \rangle \langle \mathbf{t} \rangle^T}_{=0} \\ &= \eta \langle \mathbf{x} \mathbf{t}^T - \mathbf{x} \langle \mathbf{t} \rangle^T - \langle \mathbf{x} \rangle \mathbf{t}^T + \langle \mathbf{x} \rangle \langle \mathbf{t} \rangle^T \rangle \\ &= \eta \langle (\mathbf{x} - \langle \mathbf{x} \rangle) (\mathbf{t} - \langle \mathbf{t} \rangle)^T \rangle \\ &= \langle \Delta_{cov} \mathbf{W} \rangle \end{aligned}$$

- ▶ The first term of Hebbian descent for the L2 loss includes Hebb's rule and covariance rule as special cases.
- ▶ The second term performs a normalization, lacking in Hebb's and covariance rule.

Limitation of Hebb's and covariance rule

target
output

$$\{(0, 1), (1, 0), (1, 1), (1, 1)\}$$

$$\{(0, 0), (0, 0), (1, 1), (1, 1)\}$$


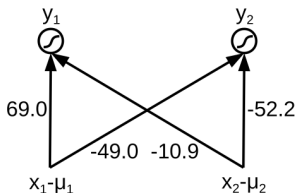
input

$$\{(0, 1), (1, 1), (1, 0), (1, 0)\}$$

(a) Hebb-rule / Covariance rule

target
output

$$\{(0, 1), (1, 0), (1, 1), (1, 1)\}$$

$$\{(0, 1), (1, 0), (1, 1), (1, 1)\}$$


input

$$\{(0, 1), (1, 1), (1, 0), (1, 0)\}$$

(b) Hebbian-descent

Sigmoid units, one update per pattern, $\eta = 10$, $b = 0$, W normalized to 100, output values rounded to two digits after the point.

- ▶ Hebb's and covariance rule do not deal well with ...
 - ▶ ... correlated patterns and
 - ▶ ... repeated patterns.

Data sets

- ▶ MNIST (LeCun et al., 1998) : 70,000 gray scale images of handwritten digits 0 to 9 with 28×28 pixels and pixel values in $[0, 1]$.
- ▶ CONNECT (Larochelle et al., 2010) : 67,587 binary 126-dimensional game-state patterns from the game Connect-4.
- ▶ ADULT (Larochelle et al., 2010) : 32,561 binary 123-dimensional patterns of census data to predict whether a person's income exceeds 50,000 dollar per year.
- ▶ CIFAR (Krizhevsky, 2009) : 60,000 color images of various objects with 32×32 pixels, converted to gray scale and rescaled to $[0, 1]$.
- ▶ RAND / RANDN : Random patterns with a size of 200 pixels. RAND has binary values with probabilities 0.5. RANDN has normally distributed values within $[0,1]$ with mean 0.5 and standard deviation 0.1.

Experiments

Data sets

- ▶ RAND → ADULT
- ▶ RANDN → CONNECT
- ▶ ADULT → MNIST
- ▶ CONNECT → CIFAR
- ▶ MNIST → CONNECT
- ▶ CIFAR → CONNECT
- ▶ ...

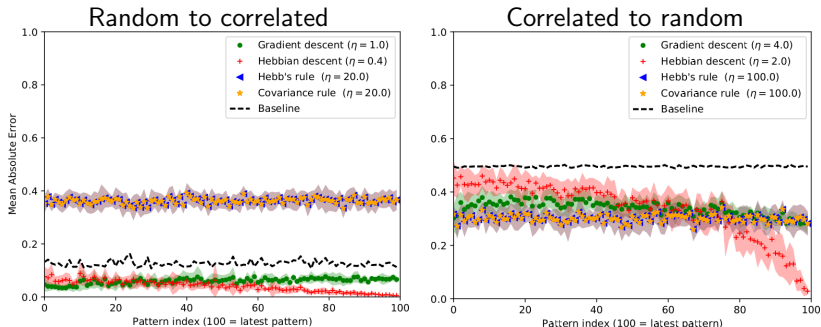
with the L2 loss.

with

Activation functions

- ▶ Linear
- ▶ Exponential Linear
- ▶ Rectifier
- ▶ Sigmoid
- ▶ Step

Associating random and correlated patterns



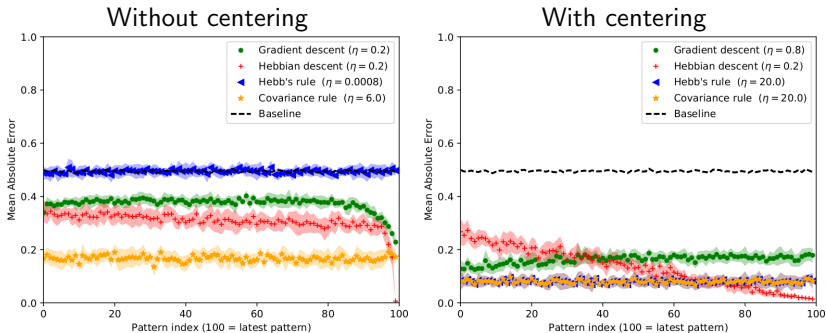
100 binary random patterns associated with 100 patterns from the ADULT dataset (binary census data) (left) and vice versa (right), one pair at a time, each pair trained once, with centering, parameters optimized for last 20 patterns.

- ▶ It is generally much easier to map random onto correlated patterns than the other way around.
- ▶ Hebb's rule and covariance rule fail due to correlations in the output.

Centering

Centering sets $\mu = \langle \mathbf{x} \rangle$ in $\mathbf{h} = \phi(\mathbf{a}) = \phi(\mathbf{W}^T(\mathbf{x} - \mu) + \mathbf{b})$

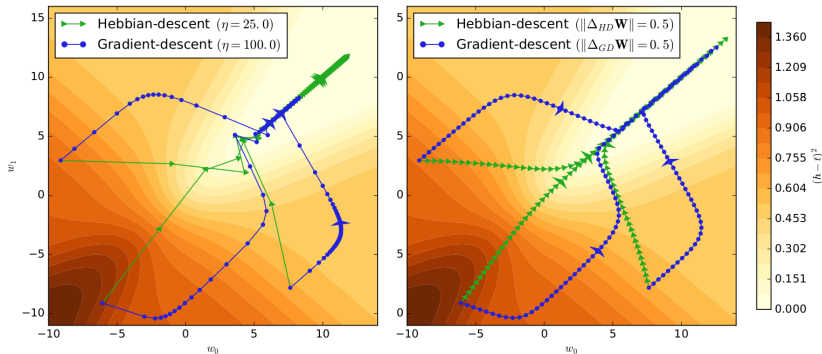
(LeCun et al, 1998; Montavon & Müller, 2012; Melchior et al., 2016)



100 binary random patterns associated with 100 binary random patterns, one pair at a time, each pair trained once.

- ▶ Centering removes the mean. Bias and weights are more independent.
- ▶ With centering Hebb's rule and covariance rule become equivalent.
- ▶ Covariance rule does not do centering since the mean is subtracted only for learning.

Hebbian vs gradient descent

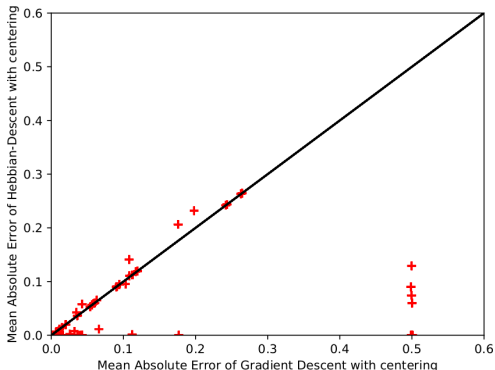


Two input units, one output unit, trained on $\{(1, 0), (0, 1), (1, 1)\} \rightarrow \{(0), (0), (1)\}$ with L2 loss for 50 full-batch updates (wider arrowheads at 25 updates) with hand tuned learning rates avoiding oscillations. Natural step size (left), constant step size (right). Equivalent results with online learning.

- ▶ In this example, Hebbian descent finds a more direct path and converges faster.

Hebbian vs gradient descent (multi epoch)

Hebbian vs gradient descent



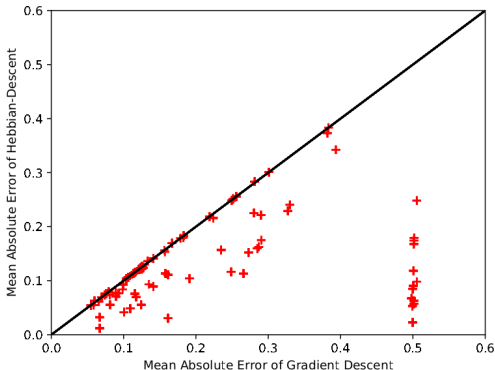
Multi epoch (with batch size one) hetero-association of 100 patterns of one data set onto 100 patterns of another data set with centering, parameters optimized for last 20 patterns on each experiment individually, 10 repetitions per experiment, mean absolute error on last 20 patterns.

For multi epoch learning:

- ▶ Hebbian descent and gradient descent perform equally well in most experiments.
- ▶ Hebbian descent has an advantage for binary target patterns.
- ▶ Gradient descent does not converge at all for the step function (values at 0.5).

Hebbian vs gradient descent (online)

Hebbian vs gradient descent

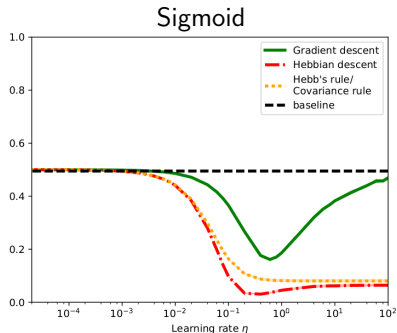
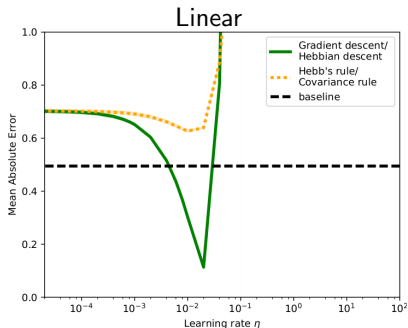


Online hetero-association of 100 patterns of one data set onto 100 patterns of another data set with centering, parameters optimized for last 20 patterns on each experiment individually, 10 repetitions per experiment, mean absolute error on last 20 patterns.

For online learning:

- ▶ Hebbian descent performs better than gradient descent in many experiments.
- ▶ The derivative of the activation function makes gradient descent inflexible.
- ▶ Gradient descent does not converge at all for the step function (values at 0.5).

Saturating activation functions are important



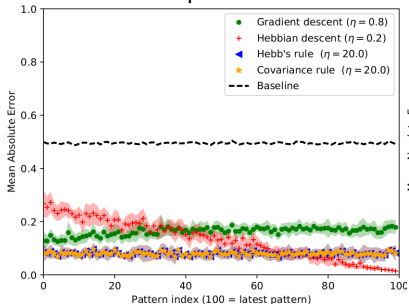
100 binary random patterns associated with 100 binary random patterns, one pair at a time, each pair trained once, with centering, parameters optimized for the last 20 patterns.

Nonlinear activation functions ...

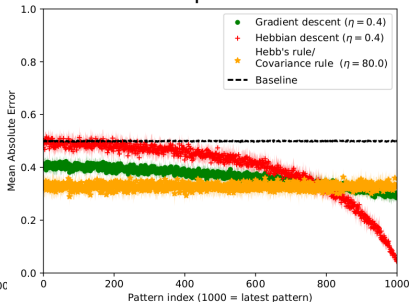
- ▶ ... adapt the output range to the target values, thereby simplifying the learning problem.
- ▶ ... make the system robust against overshooting.
- ▶ ... allow to use a wide range of large learning rates.

Gradual forgetting

100 patterns



1000 patterns

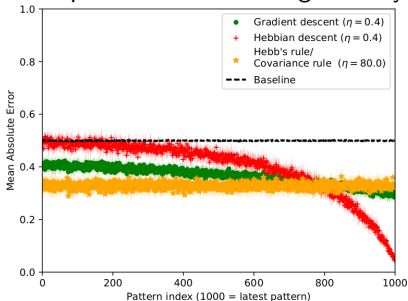


100 (left) or 1000 (right) binary random patterns associated with 100/1000 binary random patterns, one pair at a time, each pair trained once, with centering, parameters optimized for last 100 (left) or 20 (right) patterns (very similar results for last 1, 20 or all 100 patterns).

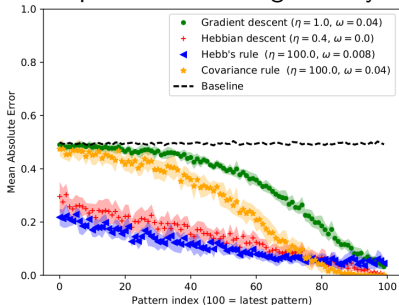
- ▶ If many more patterns are trained on than the network can store, ...
 - ▶ ... Hebbian descent gradually forgets old patterns and stores recent patterns well.
 - ▶ ... the other rules just store all patterns more or less poorly = catastrophic interference.

Weight decay

1000 patterns without weight decay



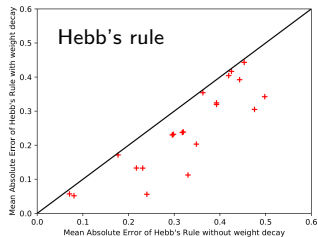
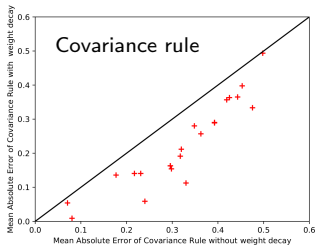
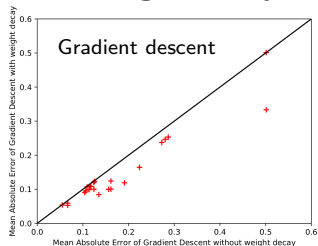
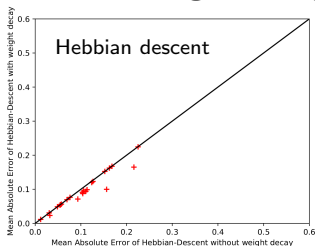
100 patterns with weight decay



1000 (left) and 100 (right) binary random patterns associated with 1000/100 binary random patterns, one pair at a time, each pair trained once, with centering, parameters optimized for last 20 patterns.

- ▶ All but Hebbian descent (optimal $\omega = 0.0$) profit from weight decay for the most recent patterns.

With weight decay vs without weight decay



Various experiments with different data sets and different activation functions, performance optimized for and measured on last 20 patterns.

- ▶ Hebbian descent profits least from weight decay but has overall best performance.
- ▶ Presumably Hebbian descent forgets more specifically than weight decay.
- ▶ Thus, weight decay is no alternative to Hebbian descent for gradual forgetting.

Summary

For hetero-association we have shown analytically:

- 1. Hebbian-descent (HD) is gradient descent with dropped $\phi'(a)$.**
- 2. HD equals gradient descent on a different loss function, if $\phi'(a) > 0$, thus**
- 3. HD provably converges for $\phi'(a) > 0$.**
- 4. HD for L2 loss with sigmoid leads to cross entropy loss. HD generalizes this idea.**

Furthermore, for L2 loss we have shown:

- 5. HD is the difference of a supervised and an unsupervised Hebb/covariance rule.**
- 6. HD can also be viewed as contrastive learning where the hidden units are clamped to the same values in positive and negative phase.**
- 7. HD optimizes a generalized linear model, if $\phi(a)$ is invertible and integrable, thus**
- 8. HD loss can be seen as the general log-likelihood loss.**

Summary

Our empirical results suggests that:

9. Centering always helps.
10. HD outperforms Hebb's/covariance rule in general.
11. HD performs like gradient decent in multi-epoch learning.
12. HD performs clearly best in online/one-shot learning.
13. HD even works with the step function, for which $\phi'(a) = 0$ for $a \neq 0$.
14. HD shows an inherent curve of forgetting, better than weight decay.

For auto-associative learning we have shown:

15. HD is related to contrastive learning in restricted Boltzmann machines.
16. HD corresponds to the encoder related part of gradient descent.
17. HD can be understood as non-linear Oja's/Sanger's rule.
18. HD does not have any objective function, but still it converges.
19. HD learns a hidden representation with equal average activity across units.
20. HD performs slightly worse than gradient descent but is computationally cheaper and biologically more plausible.

Thank you!