# Blackboard 5-1a: Saddle Points

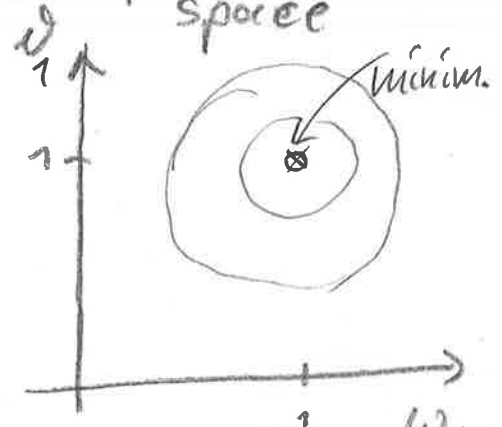**1 neuron :** input space (data)     parameter space



$\vec{w}_1 = (\omega_{11}, \omega_{12})$

$= (1, 0)$

$\vartheta_1 = 1$

normalized

$\|\vec{w}\| = 1$

$\vartheta = 1$

$\vec{w}_1$

minim.

contour lines $\omega_{11}$
of error f.
projection to $\omega_{12} = 0$

---

**2 neurons**



$\vec{w}_1 = (1, 0) ; \vartheta_1 = 1$

$\vec{w}_2 = (0, 1) ; \vartheta_2 = 1$

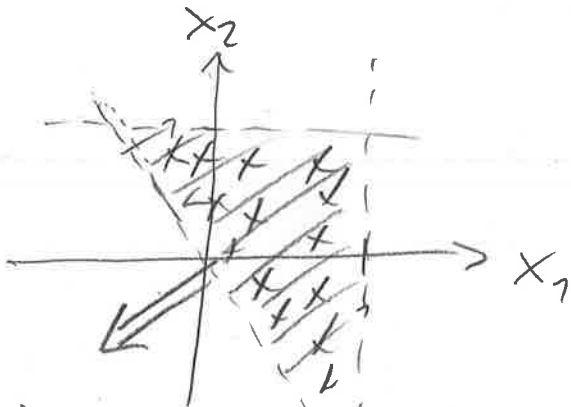$\omega_{21}$     E     $\omega_{11}$

permutation of weight vectors:
$\vec{w}_1 \rightleftarrows \vec{w}_2$

"2nd neuron implements
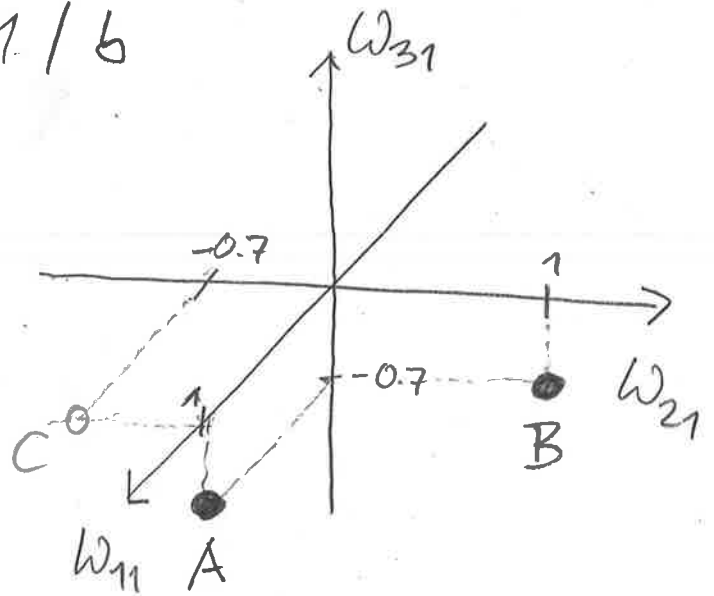first hyperplane and viceversa"

Saddle
between
minima

Blackboard 5.1/b



$\vec{W}_3 = \frac{1}{\sqrt{2}}(-1,-1)$

$\vartheta_3 = 0$

Attention: parameters
$W_{12}, W_{22}, W_{32}, \vartheta_1, \vartheta_2, \vartheta_3$
are not plotted.
Idea: adjusted
to minimum
given $W_{11}, W_{21}, W_{31}$

$$\begin{array}{ccc} W_{11} & W_{21} & W_{31} \end{array}$$
$$A = (1, 0, -0.7)$$

first permutation
$$\vec{W}_1 \rightleftharpoons \vec{W}_2$$

$$B = (0, 1, -0.7)$$

Second mutation, starting
from A: $\vec{W}_2 \rightleftharpoons \vec{W}_3$

$$C = (1, -0.7, 0)$$

etc.

saddle points between all
permutations!

Blackboard 5.2 : Momentum $E$



---- no momentum
— $d = 0.3$
$\vec{\omega}(1)$

$\omega_2 - \omega_2^*$

$\omega_2^* \quad \omega_2$

$\omega_1 - \omega_1^*$

$\dfrac{\partial E}{\partial \omega_1} = c$

$\omega_1^* \qquad \omega_1$

$$\Delta \omega_1(1) = -\gamma \cdot \frac{\partial E}{\partial \omega_1} = -\gamma c$$

$$\Delta \omega_1(2) = -\gamma c + d \cdot \Delta \omega_1(1) = -\gamma c(1+d)$$

$$\Delta \omega_1(3) = -\gamma c + d \cdot \Delta \omega_1(2) = -\gamma c - d \cdot \gamma c(1+d)$$
$$= -\gamma c\,[1 + d + d^2]$$

$$\Delta \omega_1(n) = -\gamma c\,[1 + d + \dots + d^{n-1}]$$

$$\Delta \omega_1(\infty) = -\frac{\gamma}{1-d} \cdot c$$
$$= -\gamma_{eff}\,\frac{\partial E}{\partial \omega_1}$$

# Blackboard 5.3 : stochastic gradients

projections 2D



$$\langle \Delta \omega_1 \rangle = 1 \quad ; \sqrt{\langle \Delta \omega_1^2 \rangle} \approx 1.01 \qquad \text{ratio} = 1$$

$$\langle \Delta \omega_2 \rangle = 0.1 \; ; \sqrt{\langle \Delta \omega_2^2 \rangle} = 0.1 \qquad = 1$$

$$\langle \Delta \omega_3 \rangle = 0.05 \; ; \sqrt{\langle \Delta \omega_3^2 \rangle} = \sqrt{0.5} \approx 0.7 \qquad \approx 0.07$$

$\Rightarrow$ smaller steps in "noisy" directions

note: absolute size of gradient irrelevant