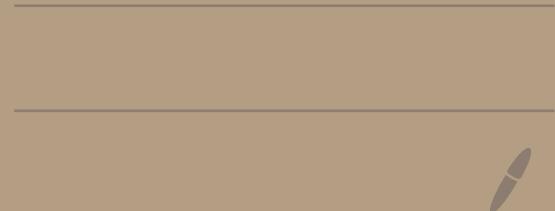
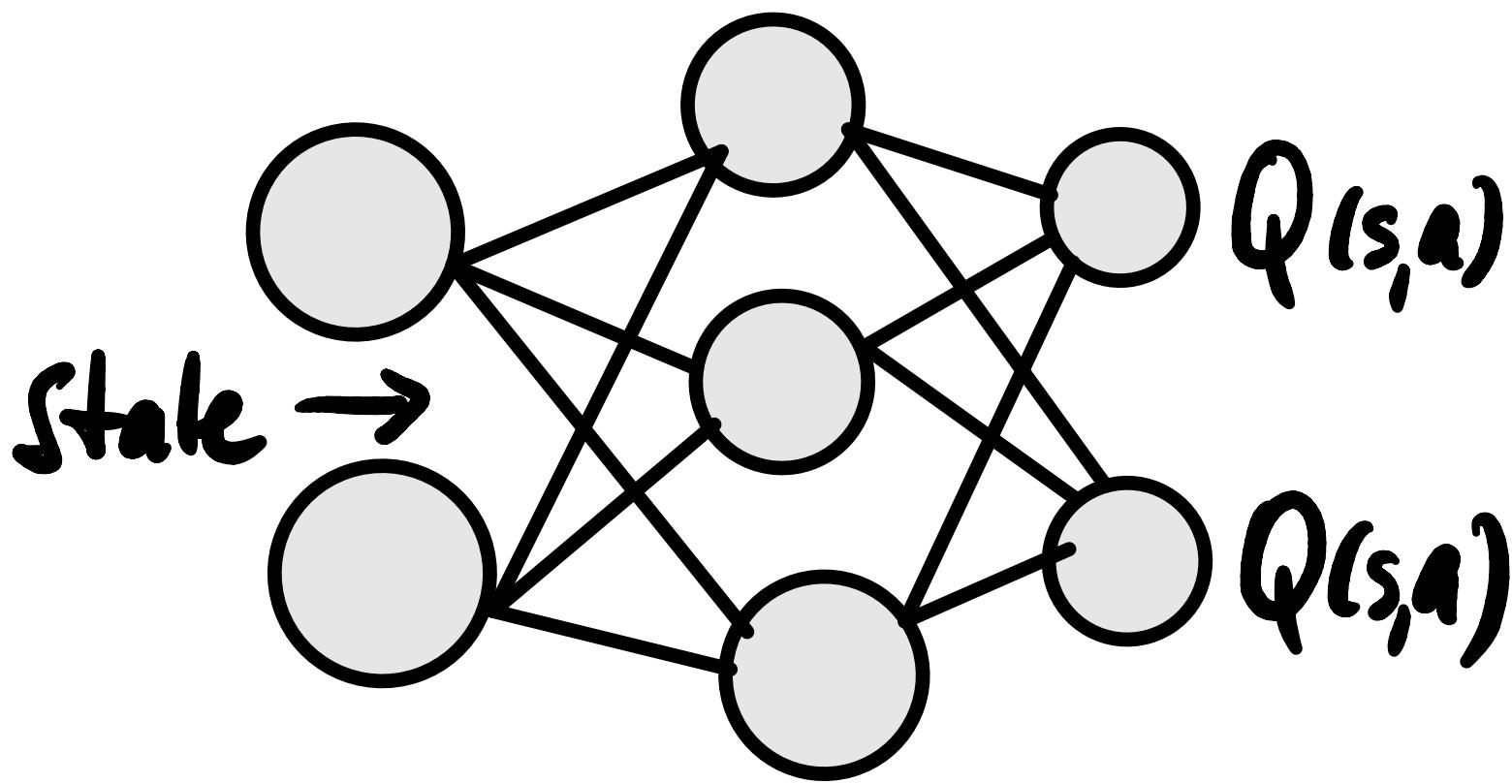


REINFORCE

09.02.2021



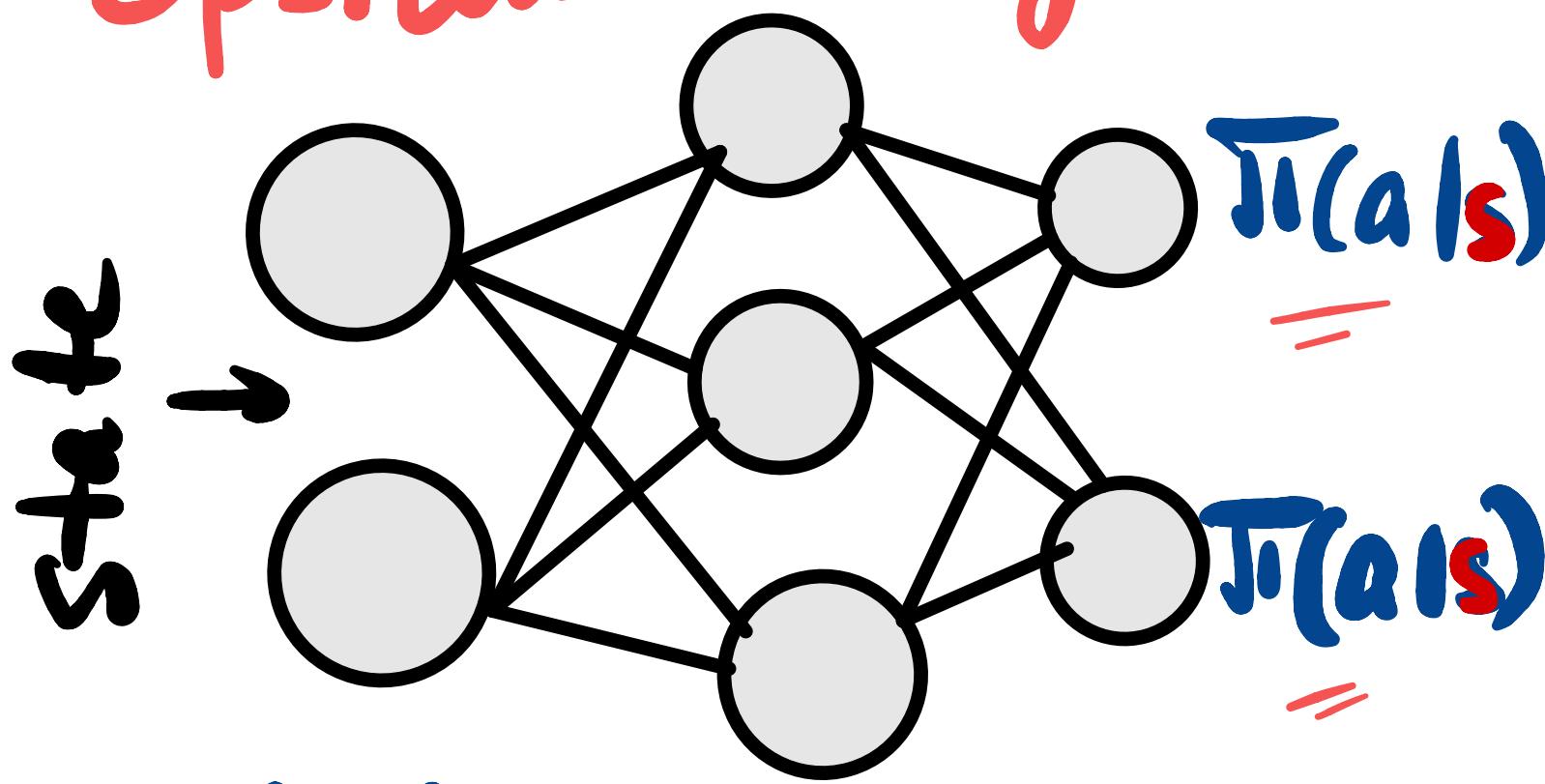
Value - Based Methods



$$a^* = \arg \max_a Q(s,a)$$

Policy: action that gives the best value

Policy Based Epsilon greedy ?

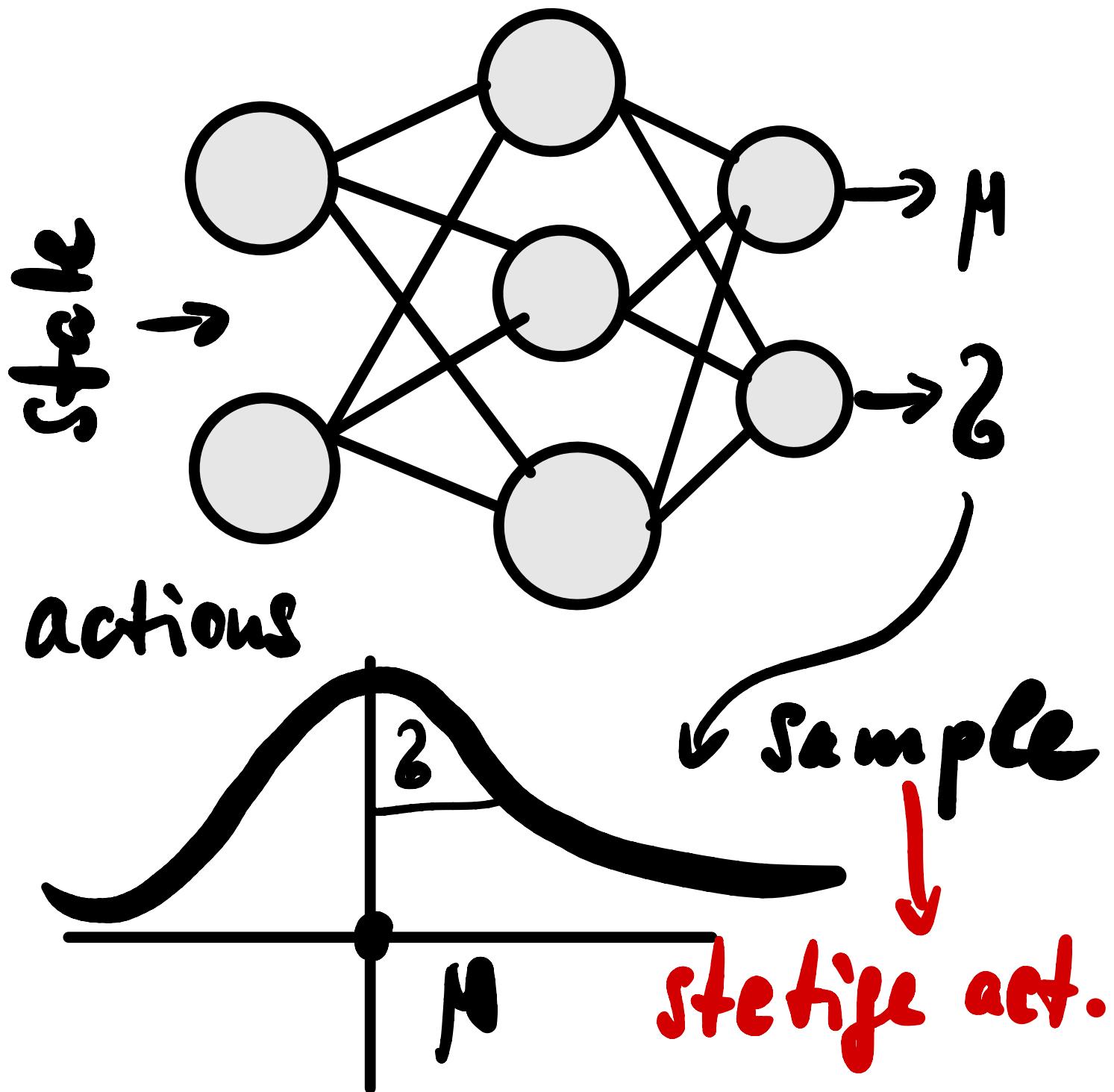


Probability we should
perform the action
given the input state

DQN - discrete
action space
each output -
action

Policy Networks -
can output
parameters of
any distribution
e.g. μ and σ

Policy



Training a Policy

DQN:

$$E \left[(G - Q(s, a))^2 \right]$$

$\xrightarrow{\rightarrow 0}$

dqn output
return

and

Kommt aus NN mit was vergleid.

$\pi_{\theta}(a | s)$?

Q - Values bzw.

$$G = R_0 + \gamma^1 \cdot R_1 + \dots \gamma^n R_n$$

aus dem Spiel



θ - Gewichte

des NNs

$$\min \left(\underbrace{f(G) - \bar{\pi}}_{?} \right)^2$$

Policy Gradient Derivation

1. **max**

$$\sum_i p_i \cdot R_i$$

$G = E[R_1 + R_2 + \dots + R_T]$
 $=$ *return \rightarrow exp. value*

2. state-action pairs

$$\Omega = \{(s_1, a_1), \dots, (s_T, a_T)\}$$

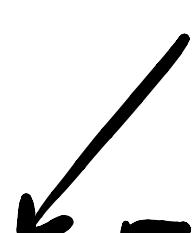
3. $G(\Omega) = \sum_t R(s_t, a_t)$
return

Gradient

ASCENT

$$\theta_{\text{new}} = \theta + \nabla J_{\theta} \cdot LR$$

ascent



LR - Learning rate

$$\nabla J_{\theta} \rightarrow$$

$$\theta_{\text{new}} = \theta_{\text{old}} - (\nabla J_{\theta} \cdot L)$$

4. Objective $J(\theta)$

$$J(\theta) = \sum_{\Omega} P(\Omega; \theta) \cdot G(\Omega)$$

Diagram illustrating the components of the objective function:

- Network params** (blue) point to the θ in the equation.
- Summe über State-actions** (red) points to the summation symbol (\sum) and the Ω below it.
- Probs** (blue) point to the $P(\Omega; \theta)$ term.
- Policy which want to find** (red) points to the $G(\Omega)$ term.
- Return** (red) points to the $G(\Omega)$ term.

Summe über

State - actions

We want to maximize our objective Funktion
maximize expected return!

5. Objective Gradient

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \sum_{\Omega} P(\Omega; \theta) \cdot G(\Omega)$$

$\frac{\partial}{\partial \theta} (\Omega)$

$$\nabla_{\theta} J(\theta) = \sum_{\Omega} \frac{\nabla_{\theta} P(\Omega; \theta) \cdot G(\Omega)}{1}$$

$$(a+b)' = a' + b'$$

$$\nabla(a+b) = \nabla a + \nabla b$$

6. Trick :

$$\nabla_{\theta} J(\theta) = \sum_{\Omega} P(\Omega; \theta) \cdot$$

$$\frac{\nabla P(\Omega; \theta)}{P(\Omega; \theta)} \cdot$$

$$G(\Omega)$$

$$\nabla \log f = \frac{\nabla f}{f}$$

$$\nabla_{\theta} J(\theta) = \sum_{\Omega} \underbrace{P(\Omega; \theta)}_{\text{probs}} \nabla_{\theta} \log \cdot P(\Omega; \theta) \cdot G(\Omega)$$

Summe

Erwartungswert

$$\ln(\underline{x})' = \frac{E \underline{1}}{\underline{x}} = \frac{\checkmark f}{f}$$

$$\nabla_{\theta} J(\theta) = E \left[\nabla_{\theta} \log P(\Omega; \theta) \cdot \frac{G(\Omega)}{G(\Omega)} \right]$$

1. Gradient tries to make P more probable
2. IF G is positive, reinforce the actions, that lead to it

$\log P \cdot G \uparrow \textcircled{1}$

$\log P \downarrow \textcircled{2} \cdot G \downarrow \textcircled{1}$

$$G = p(r + \gamma \max_a Q)$$

direct

BEQ

mur dis krek
actionen

Value Iter

SARS(A)
(Tabular Q-L.)

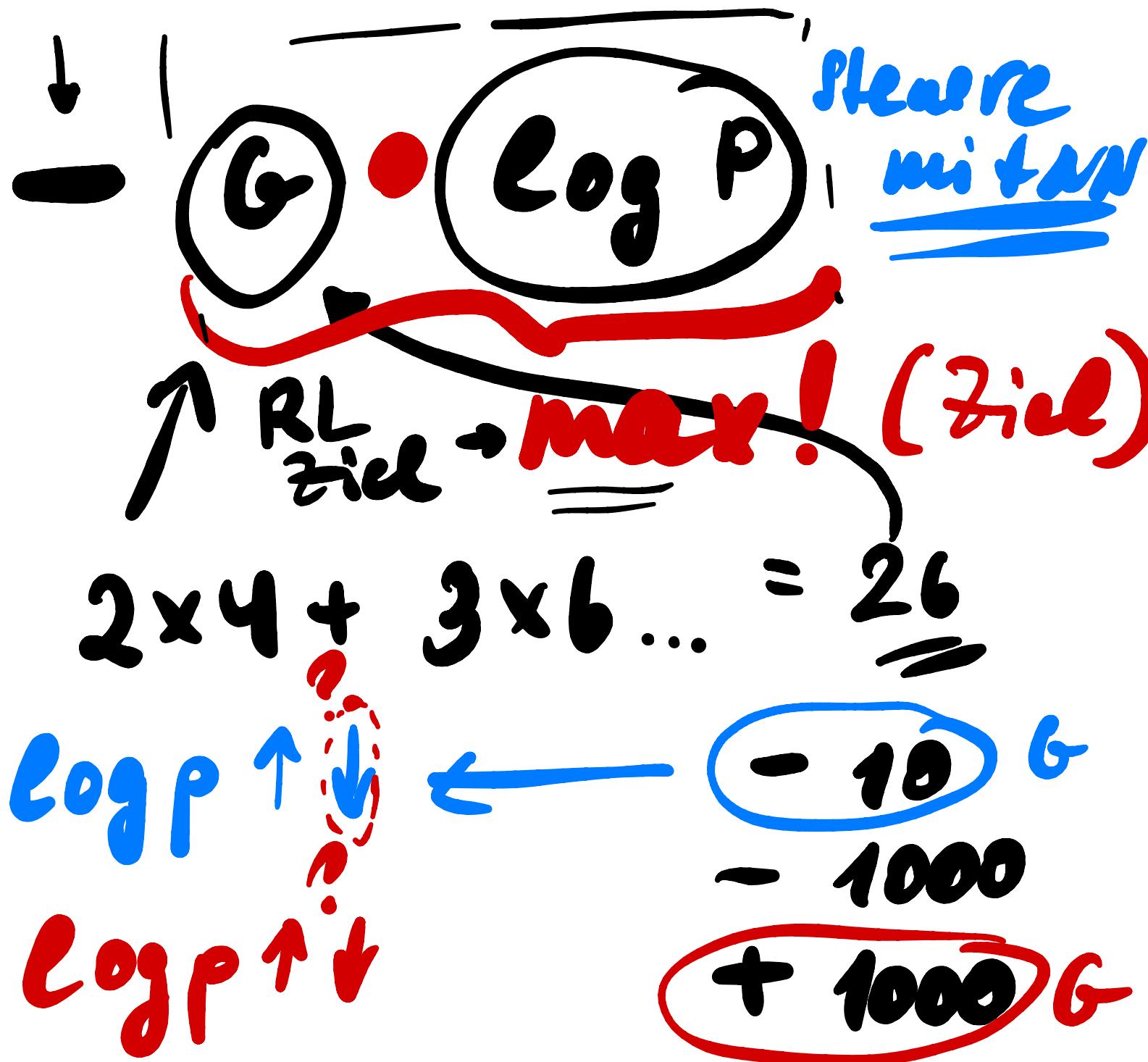
DQN } Paper
DDQN } 2013
RDQN }

Value Iter

-->
direkt policy
ausrechnen
 $\pi(a|s)$
=====

Familie von

Policy Gradient
Verfahren



NN \rightarrow Dicmer

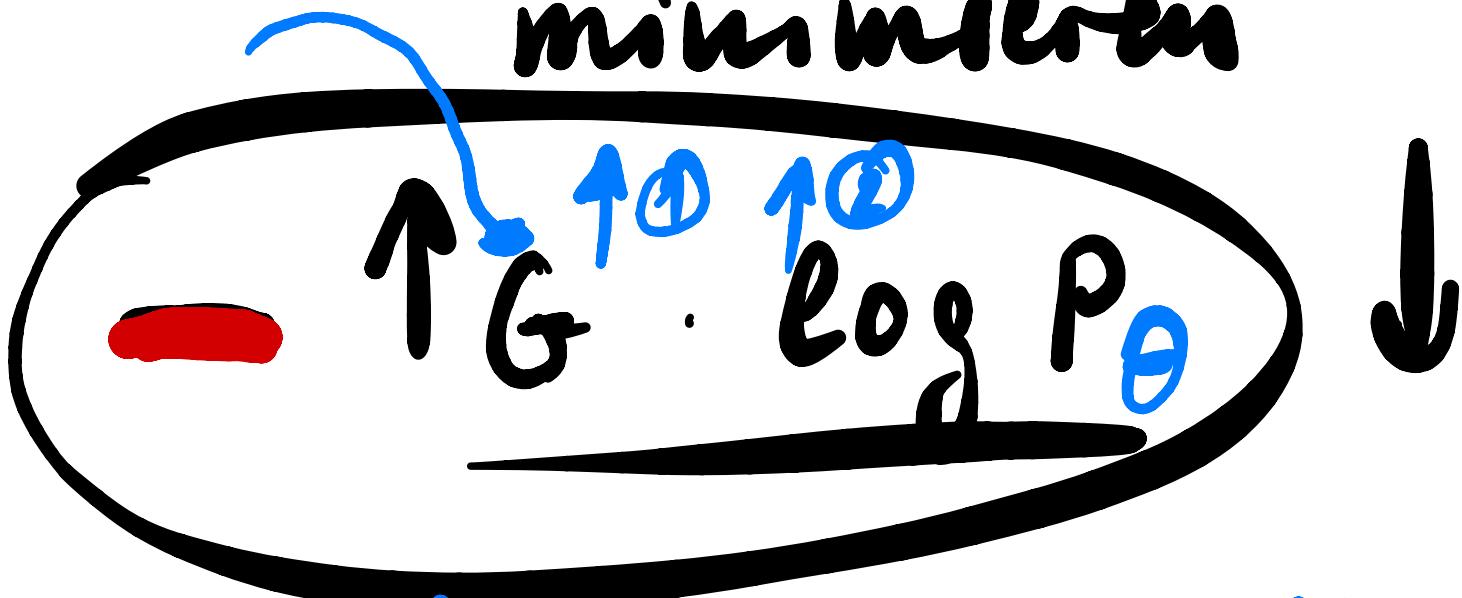
RL - wir

wir: $G \cdot \log P \uparrow$

max

NN: kann nur
minimieren

$G \cdot \log P_\theta$



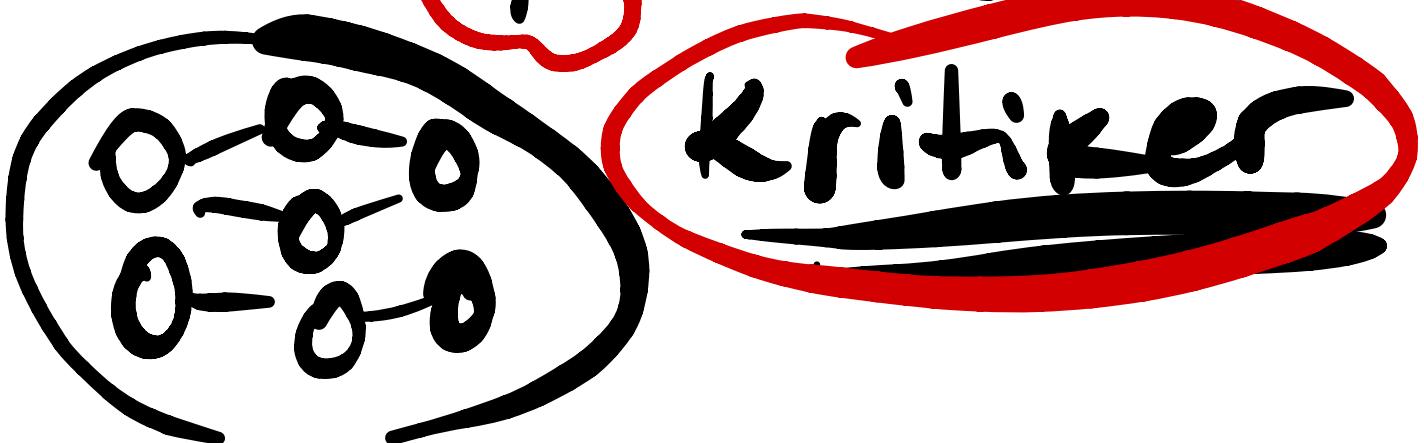
REINFORCE actions
mit $G + re$

Spoiler :)



Statis Advantage

$A = Q_{S,A}$ - V_S (Schon da)



A2C

Actor (REINF.)

Critic (PQN)

Log P. A (statt
c)

$A = Q(\text{new}) - V_{\text{last}}$
Vergleich.
von Actor)

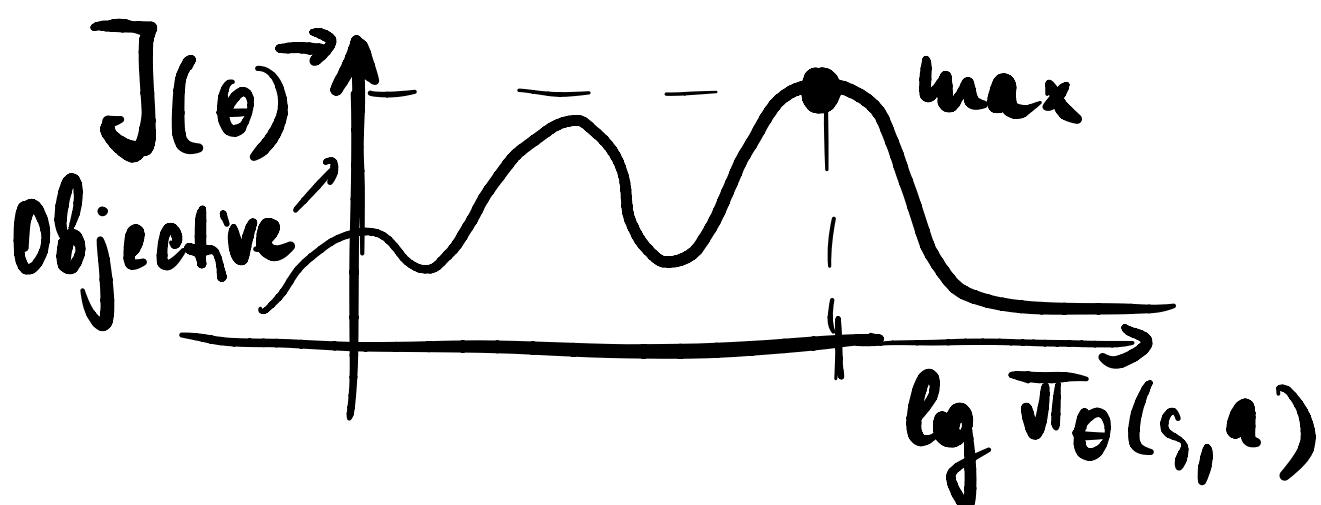
$$L(\theta) = -\log P(\Omega; \theta) \cdot G(\Omega)$$

Objective \Rightarrow loss

By setting it negative we perform Gradient ascent, looking for min of the negative Loss

$$\begin{aligned}\theta_{t+1} &= \theta_t - (-\nabla \log \pi_\theta(s, a) \cdot \hat{Q}(s, a)) \\ &= \theta_t + \nabla \log \pi_\theta(s, a) \cdot \hat{Q}(s, a)\end{aligned}$$

Policy Grad.



Kyber. Regelkreis

FB

FF

Env

$\pi_{\text{neu}} \rightarrow \pi$

π

reward

Ergebnis von π

π_{neu}

action

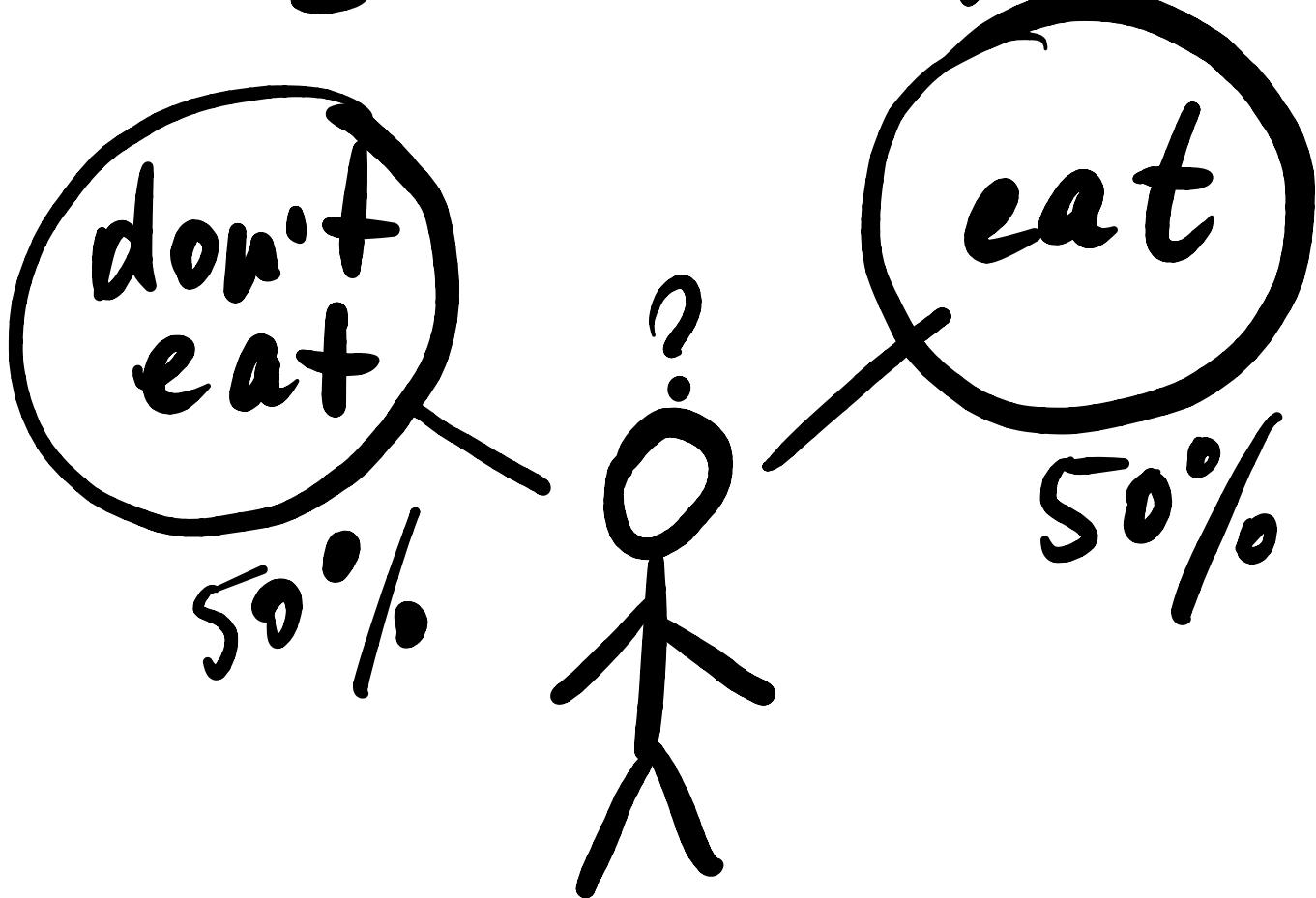
NN
 $\max \log \pi_i \cdot \hat{Q}$

$$\frac{-\log \pi \cdot \hat{Q}(s, a)}{-\log p \cdot G}$$

Loss

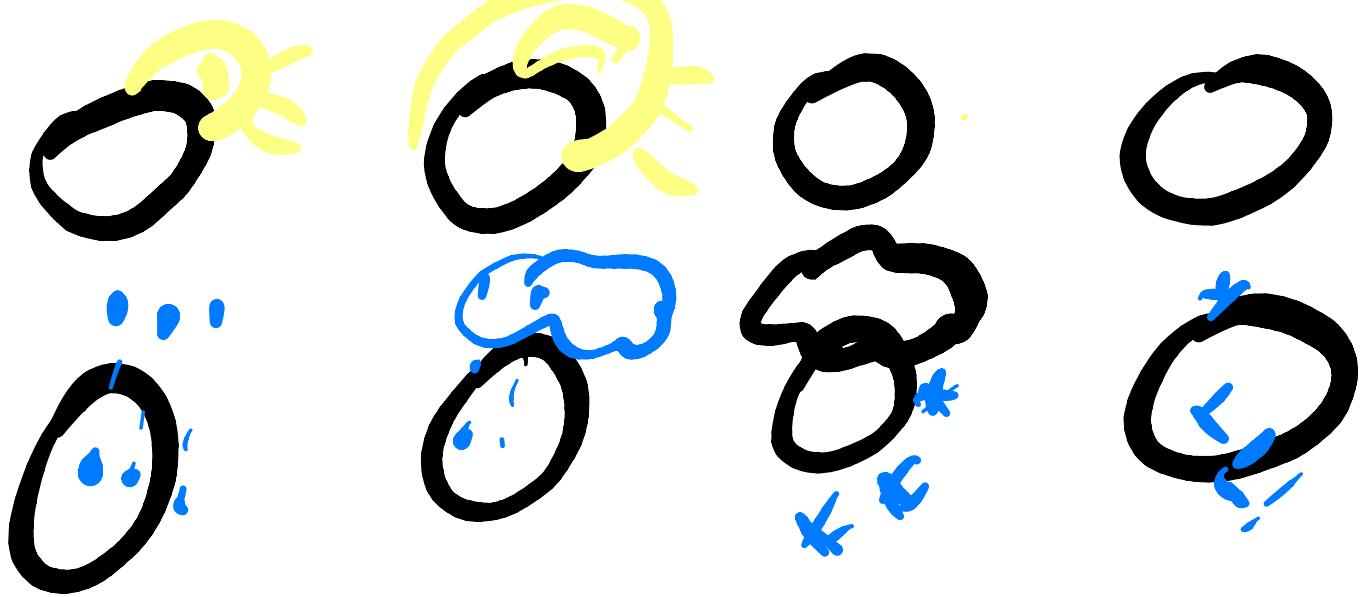
Claude Shannon

1 bit 0/1



→ 0 DE
1 E

Ursicherheit um faktor
② reduziert



1. Entropy
2. CrossEntropy
↳ Objective
classification
3. KL Divergenz
↳ Autoencoder

0000
0000

$$2^3 = 8$$

3 bit

gleich
wahrsch.

$$\log_2 8 = 3$$

Message

Length

$$-\log_2 \frac{1}{8} = \log_2 8$$

p

$\sum p_i - (\log_2 p_i)$

↓ ↓

message
Length

Wahrsch. Entropy

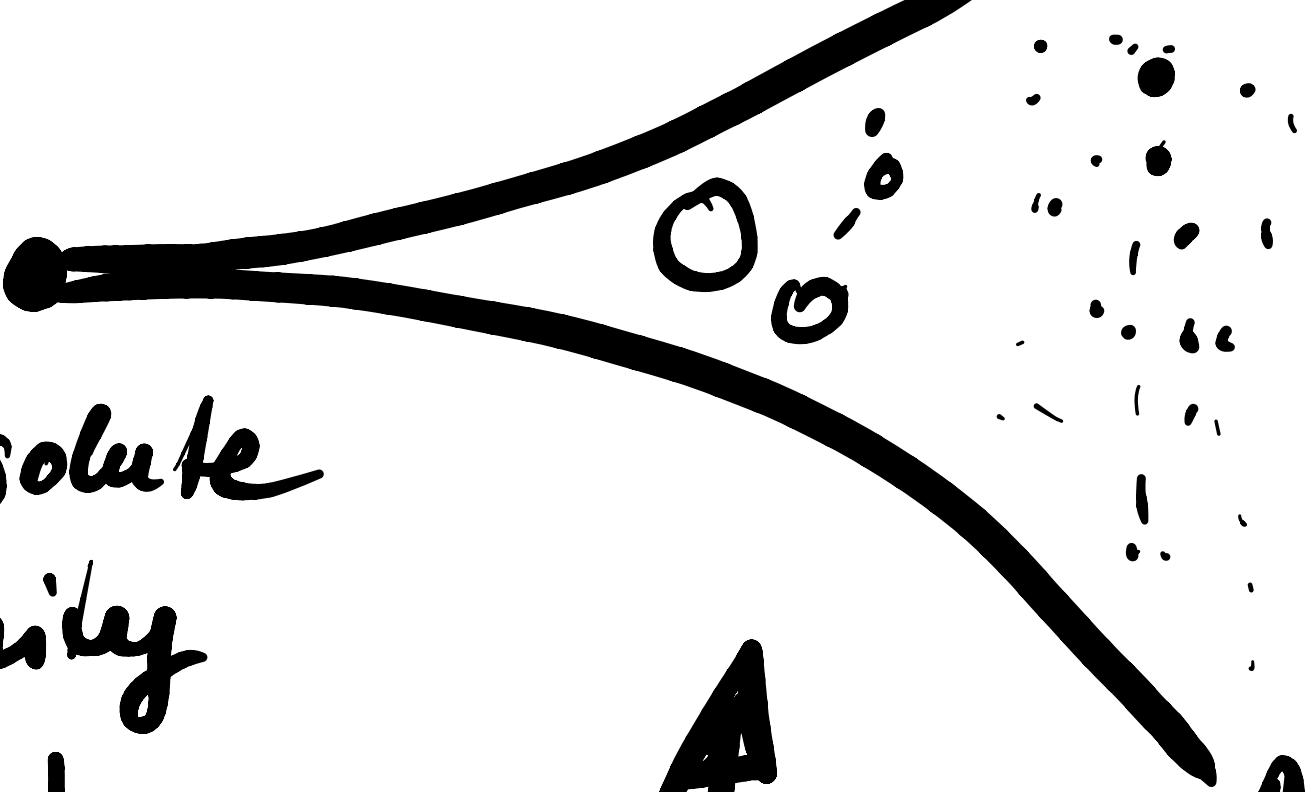
je komplexer
desto länger message
um das zu beschreiben

$$\delta \rightarrow \frac{1}{\delta}$$

$$10.000 \quad \frac{1}{10.000}$$

$$-\log \frac{1}{10.000} x + \text{gross}$$

$$\lim_{P \rightarrow 0} -\log P \rightarrow \underline{+\infty}$$



Absolute
Unity

Entropy = 0 ↑ Entropy ↑

message ↑

DNS →

0
75%

0
25%

$$-\ln 0.75 - \ln 0.25$$

message
length

message.

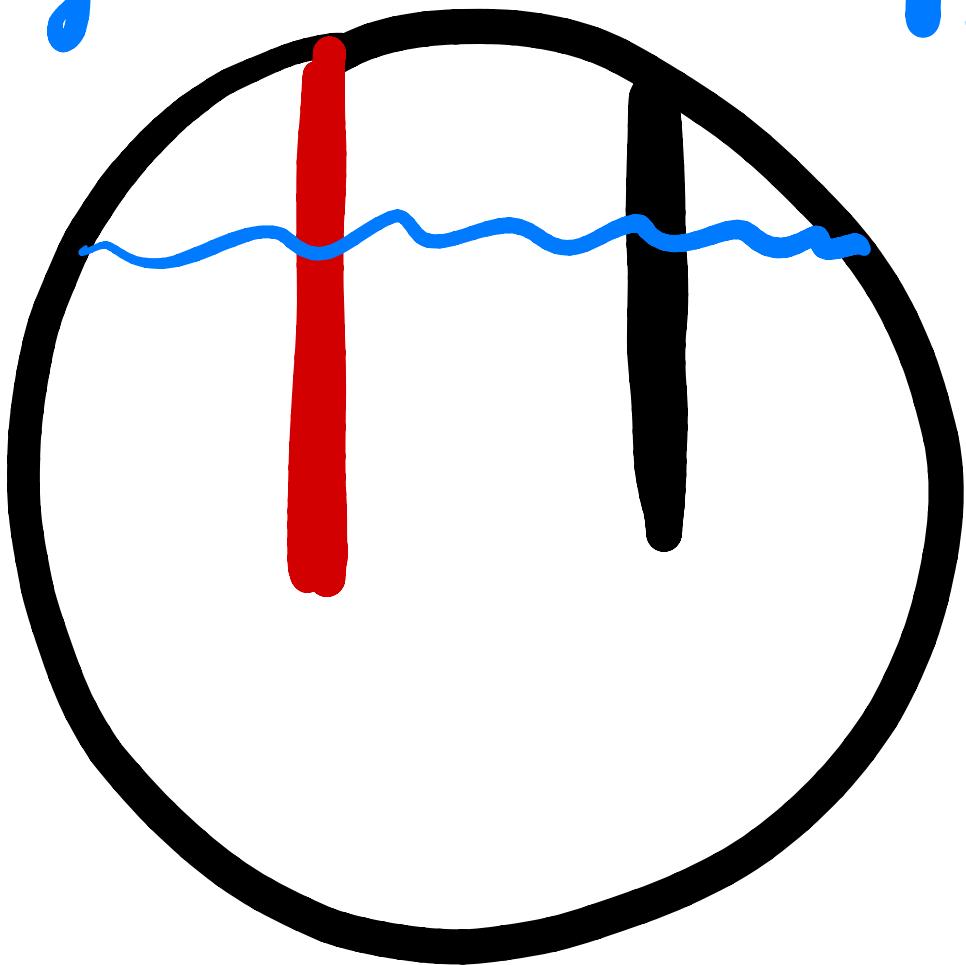
$$0.75 \cdot \ln \frac{1}{0.75} + 0.25 \ln \frac{1}{0.25}$$

Number Bits Bits

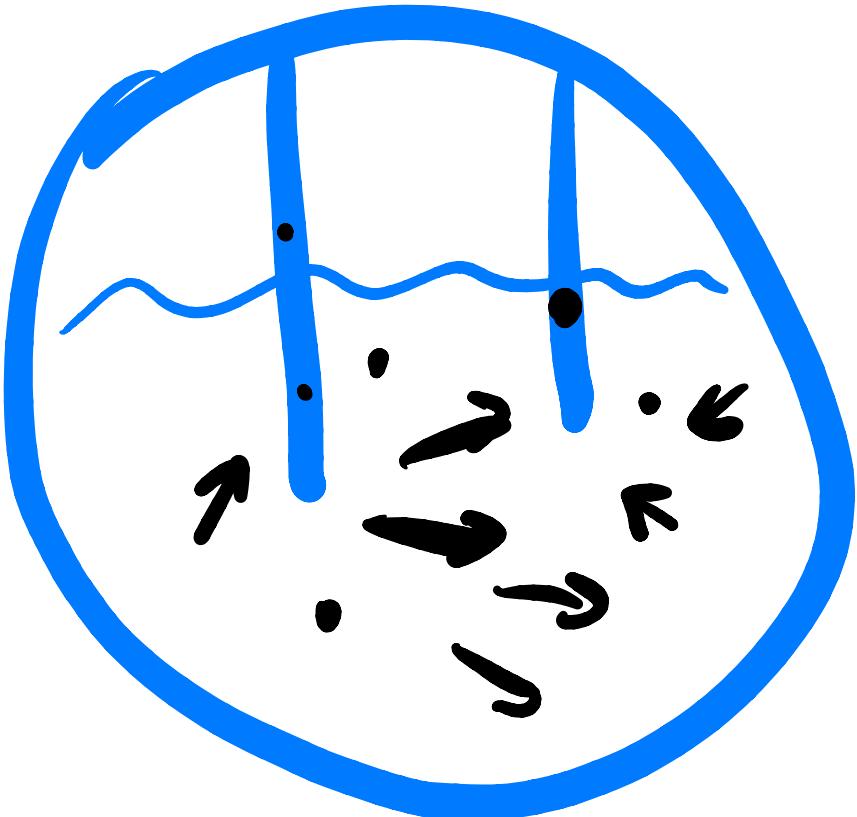
Erwartungswert Mess-
Länge

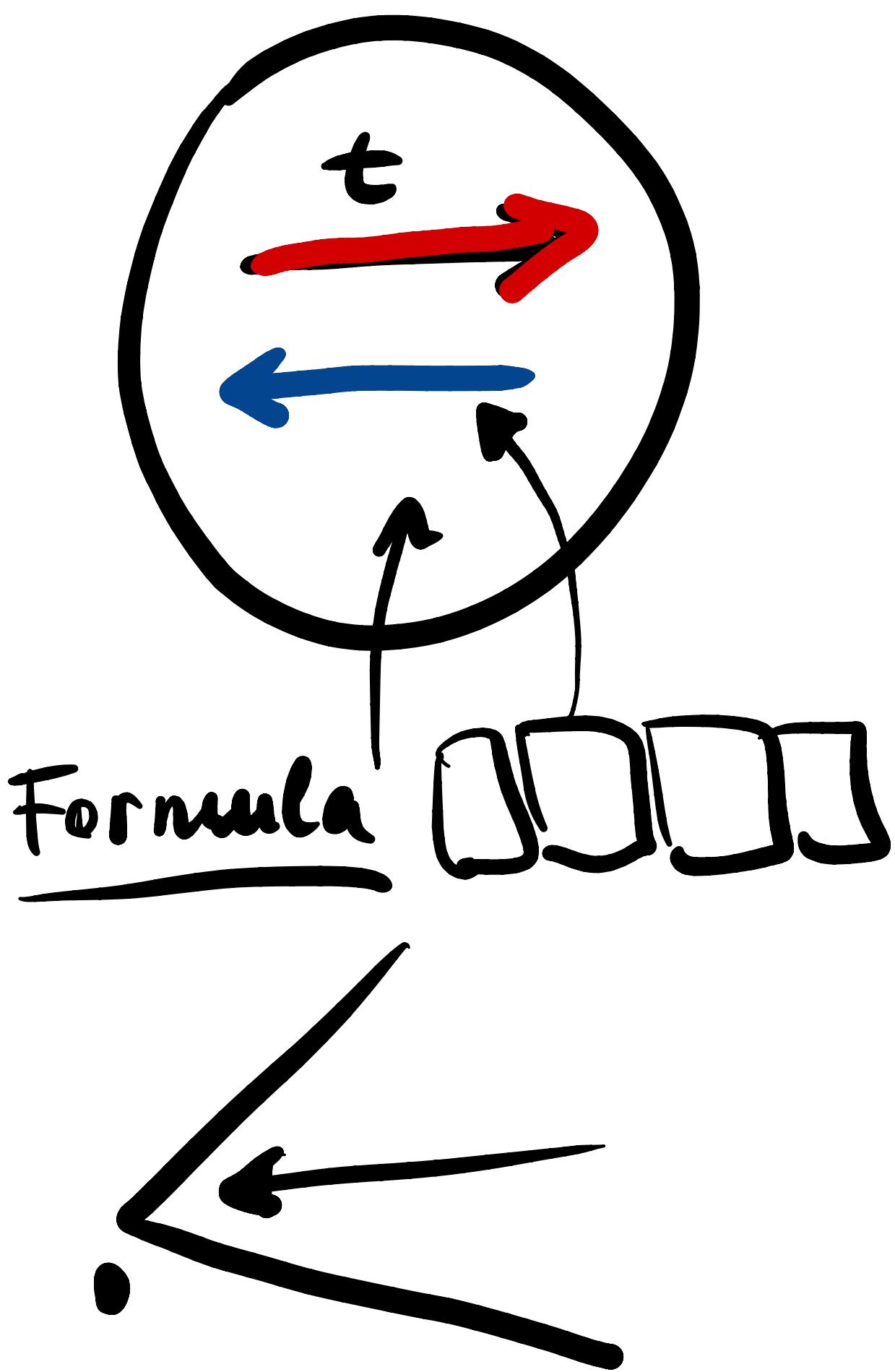
Physik

1.



2.





2. Cross Entropy

$$-\sum p \cdot \ln p$$

Wahr Komplexität

$$-\sum p \cdot \ln q$$

→ cross Entropy

$\ln q$ - unsere Schätz.
von Message 1.
- Komplexität

$H(p)$ Entropie

$H(p, q)$ Cross Entropie

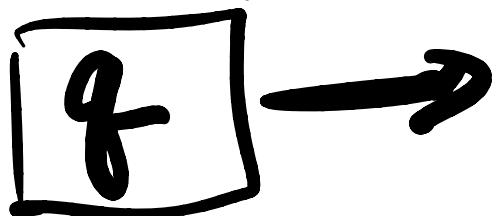
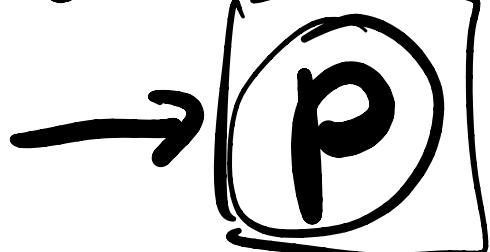
target

\hat{y}

$KL = H(p, q) - H(q)$
Variational

→ Kullback Leibler

Divergenz
Schätzungsfehler



VAE

$D_{in} \neq D_{out}$

Verteilung p $\neq q$ Verteilung
Stat. properties

$\hookrightarrow H(p, q) = H(p)$

$\hookrightarrow KL = 0.$!

Supervised

Obj. Cross Entr.



geschätzte $H(\rho, \hat{\rho}) = H(\rho)$

Kompl. $KL = 0$

0
Ea +
90%

0
DE
10%

$$-\left(0.9 \cdot \log 0.9 + 0.1 \log 0.1 \right)$$

$$-\left(0.9 \cdot \log 0.5 + 0.1 \log 0.5 \right)$$

$$\underline{E=mc^2}$$

$$H(p,q) > H(p)$$