

LECTURE 19

review so far

supervised!

...			
n	x_i	y_i	

learn $P_\theta(y|x)$
 max log likelihood
 $\log\text{-lik}(\theta) = \frac{1}{n} \sum_{i=1}^n \log P_\theta(y_i|x_i)$

θ

y
s
...
x

Neural Network
"a team of vectors"

$\frac{\partial \text{loss}}{\partial s_i} = -(y_i - s_i) = -e_i$
 $\frac{\partial \text{loss}}{\partial \theta}$

continuous y regression
 $-\log\text{-likelihood} = \text{Loss} = \frac{1}{n} \sum_{i=1}^n (y_i - s_i)^2 / 2$

categorical y (classification)

y	y_k	
↑	↑	
p	$p_1 \dots p_k \dots p_c$	
↑	↑ softmax	
s	$s_1 \dots s_k \dots s_c$	
↑	↑	
x		

category / word

$$\log\text{-likelihood} = \frac{1}{n} \sum_{i=1}^n \langle \begin{matrix} y_i \\ s_i \end{matrix} \rangle - \log z$$

$$z = \sum_{k=1}^c e^{s_k}$$

$$\text{Loss} = -\log\text{-likelihood}$$

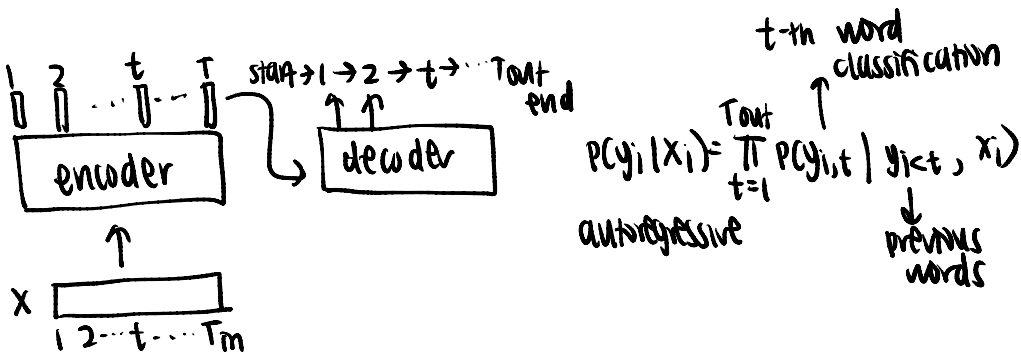
$$\frac{\partial \text{loss}}{\partial s_i} = - \left(\begin{matrix} y_i \\ p_i \end{matrix} \right) = -e_i$$

translation

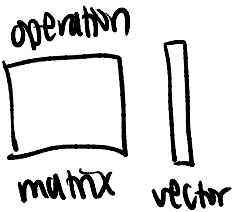
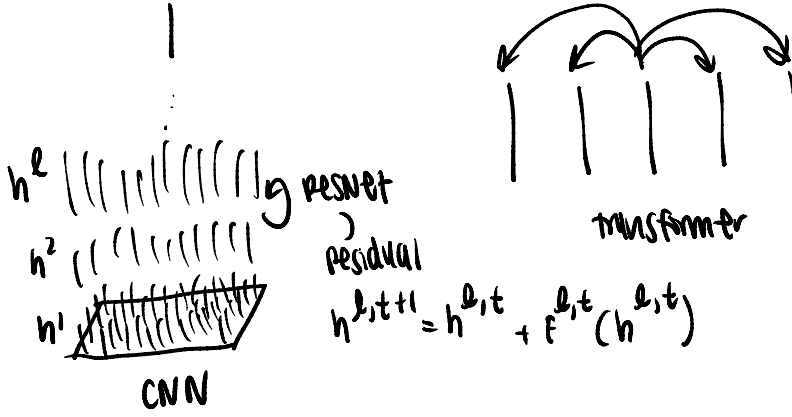
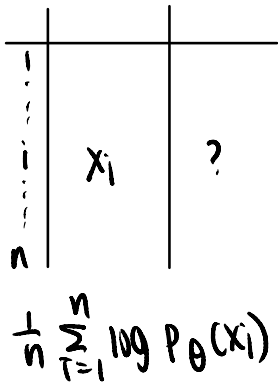
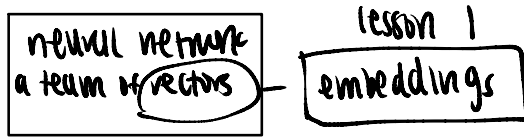
x_i	y_i
English	Spanish

spanish
English

$$\log\text{-lik}(\theta) = \frac{1}{n} \sum_{i=1}^n \log P(y_i|x_i)$$



unsupervised:



like a computer program

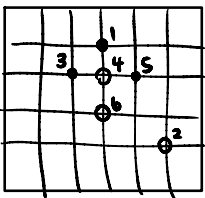


lesson 2
learned computation \rightarrow multiple steps
(instead of multiple layers)



Reinforcement learning (RL)

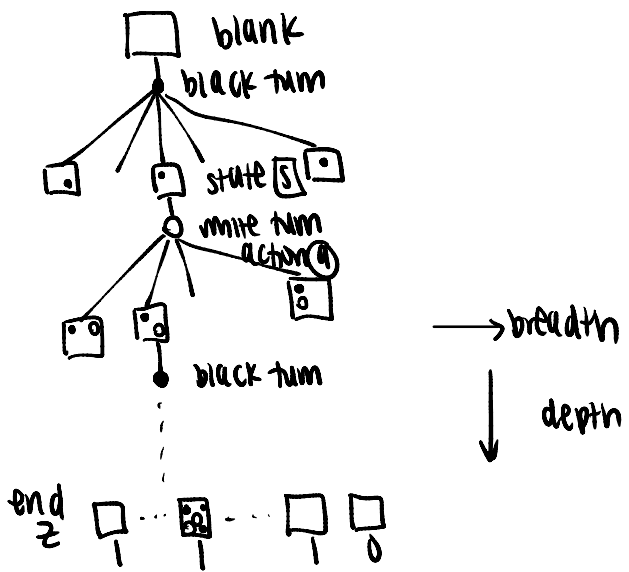
Alpha Go



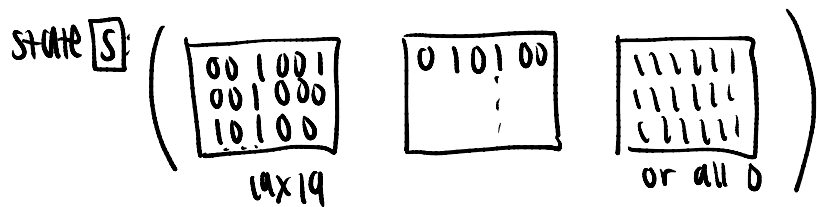
- black
- white

$$z = \begin{cases} 1 & \text{if } \bullet \text{ wins} \\ -1 & \text{if } \circ \text{ wins} \\ 0 & \text{if draw} \end{cases}$$

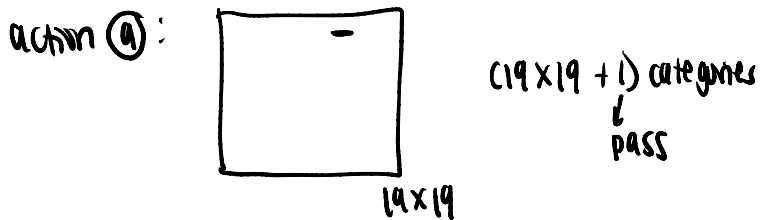
Game Tree



→ breadth
↓ depth



1: black stone 1: white stone all 1: black turn
all 0: white turn

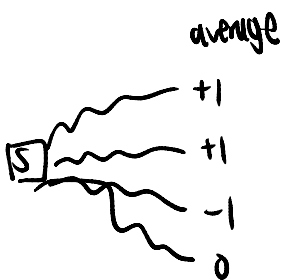
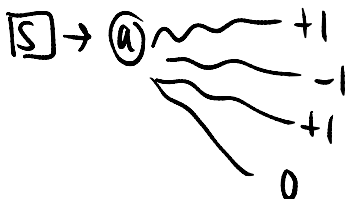


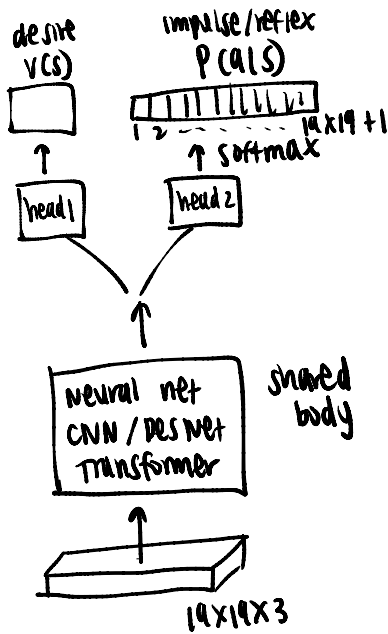
① policy: $P(A|S)$ classification

②_a value (state): $V(S) = E_{\text{policy}}(z|S)$ (e.g. .7)

$S \rightarrow A \rightarrow \dots \rightarrow \square \rightarrow z$

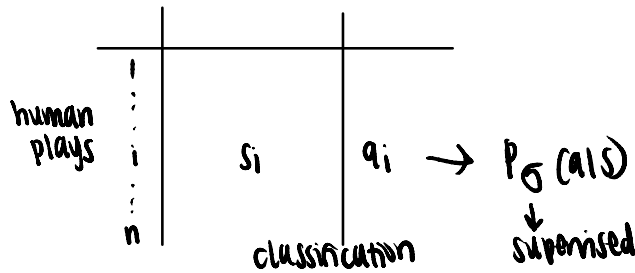
②_b value (state, action): $Q(S,A) = E_{\text{policy}}(z|S,A)$





learning:

STEP 1: supervised learning



$$\max_{\theta} \frac{1}{n} \sum_{i=1}^n \log P_{\theta}(a_i | s_i)$$

Imitation, behavior cloning
Stochastic gradient ascent

$$\Delta \theta \propto \frac{\partial}{\partial \theta} \log P_{\theta}(a_i | s_i)$$

① generated by human

STEP 2: RL

$P_{\theta}(a|s)$

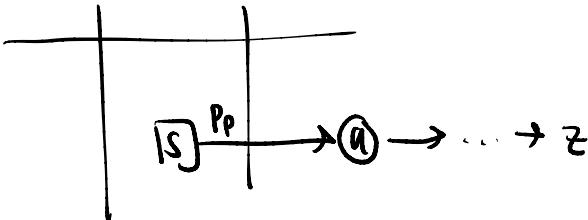
reinforcement

initialize $\theta = \theta_0$

let θ play against θ_0

repeat

each follows own policy



Policy Gradient (how we update the policy)

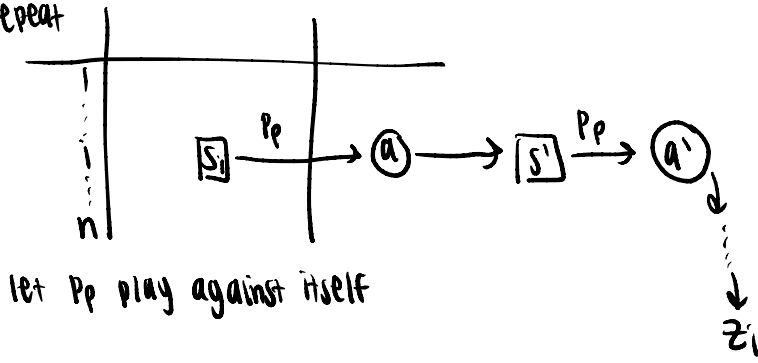
$$\Delta \theta \propto \frac{\partial}{\partial \theta} \log P_{\theta}(a|s) \cdot [z]$$

② reinforcement
 +1 → good action, ↑ prob. of action
 -1 → bad action
 0 → want ↓ prob. of action

① generated by P_{θ} self
try out different actions & see what will happen

STEP 3: value network

repeat



regression

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n (z_i - v_{\theta}(s_i))^2$$

$$\Delta \theta \propto - \frac{\partial}{\partial \theta} (z - v_{\theta}(s))^2$$