

# 10703 Deep Reinforcement Learning and Control

Russ Salakhutdinov

Machine Learning Department  
rsalakhu@cs.cmu.edu

## Policy Gradient II

# Used Materials

- **Disclaimer:** Much of the material and slides for this lecture were borrowed from Rich Sutton's RL class and David Silver's Deep RL tutorial

# Policy-Based Reinforcement Learning

- ▶ So far we approximated the value or action-value function using parameters  $\theta$  (e.g. neural networks)

$$V_{\theta}(s) \approx V^{\pi}(s)$$

$$Q_{\theta}(s, a) \approx Q^{\pi}(s, a)$$

- ▶ A policy was generated directly from the value function e.g. using  $\epsilon$ -greedy
- ▶ In this lecture we will directly parameterize the policy

$$\pi_{\theta}(s, a) = \mathbb{P}[a \mid s, \theta]$$

- ▶ We will focus again on model-free reinforcement learning

# Policy Gradient Theorem

- ▶ The policy gradient theorem generalizes the likelihood ratio approach to multi-step MDPs
  - ▶ Replaces instantaneous reward  $r$  with long-term value  $Q^\pi(s,a)$
  - ▶ Policy gradient theorem applies to start state objective, average reward and average value objective
- 

- ▶ For any differentiable policy  $\pi_\theta(s, a)$ , the policy gradient is

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) Q^{\pi_\theta}(s, a)]$$

# Monte-Carlo Policy Gradient (REINFORCE)

- ▶ Update parameters by **stochastic gradient ascent**
- ▶ Using policy gradient theorem
- ▶ Using return  $G_t$  as an unbiased sample of  $Q^{\pi_\theta}(s_t, a_t)$

$$\Delta\theta_t = \alpha G_t \nabla_{\theta} \log \pi_{\theta}(s_t, a_t)$$

## REINFORCE, A Monte-Carlo Policy-Gradient Method (episodic)

Input: a differentiable policy parameterization  $\pi(a|s, \theta), \forall a \in \mathcal{A}, s \in \mathcal{S}, \theta \in \mathbb{R}^n$

Initialize policy weights  $\theta$

Repeat forever:

    Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ , following  $\pi(\cdot|\cdot, \theta)$

    For each step of the episode  $t = 0, \dots, T - 1$ :

$G_t \leftarrow$  return from step  $t$

$\theta \leftarrow \theta + \alpha \gamma^t G_t \nabla_{\theta} \log \pi(A_t|S_t, \theta)$

# Actor-Critic

- ▶ Monte-Carlo policy gradient still has **high variance**
- ▶ We can use a **critic** to estimate the action-value function:

$$Q_w(s, a) \approx Q^{\pi_\theta}(s, a)$$

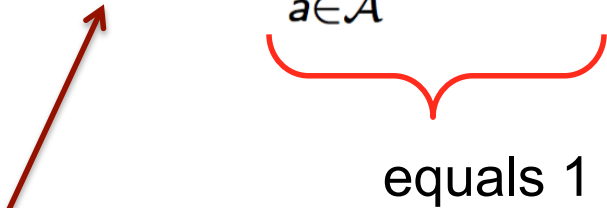
- ▶ **Actor-critic algorithms** maintain two sets of parameters
  - **Critic Updates** action-value function parameters  $w$
  - **Actor Updates** policy parameters  $\theta$ , in direction suggested by critic
- ▶ Actor-critic algorithms follow an approximate policy gradient

$$\nabla_\theta J(\theta) \approx \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)]$$

$$\Delta\theta = \alpha \nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)$$

# Reducing Variance Using a Baseline

- ▶ We can subtract a **baseline function**  $B(s)$  from the policy gradient
- ▶ This can reduce variance, **without changing expectation!**

$$\begin{aligned}\mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) B(s)] &= \sum_{s \in \mathcal{S}} d^{\pi_{\theta}}(s) \sum_a \nabla_{\theta} \pi_{\theta}(s, a) B(s) \\ &= \sum_{s \in \mathcal{S}} d^{\pi_{\theta}} B(s) \nabla_{\theta} \underbrace{\sum_{a \in \mathcal{A}} \pi_{\theta}(s, a)}_{\text{equals 1}} \\ &= 0\end{aligned}$$


Function of state  $s$ , but  
not action  $a$

- ▶ A good baseline is the state value function  $B(s) = V^{\pi_{\theta}}(s)$

# Reducing Variance Using a Baseline

- ▶ We can subtract a **baseline function**  $B(s)$  from the policy gradient
- ▶ This can reduce variance, **without changing expectation!**

$$\mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) B(s)] = 0$$

- ▶ A good baseline is the state value function  $B(s) = V^{\pi_{\theta}}(s)$
- ▶ So we can rewrite the policy gradient using **the advantage function**:

$$A^{\pi_{\theta}}(s, a) = Q^{\pi_{\theta}}(s, a) - V^{\pi_{\theta}}(s)$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) A^{\pi_{\theta}}(s, a)]$$

- ▶ Note that it is the exact same policy gradient:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q^{\pi_{\theta}}(s, a)]$$



# Estimating the Advantage Function

- ▶ The advantage function can significantly reduce variance of policy gradient
- ▶ So the critic should really estimate the advantage function
- ▶ For example, by estimating both  $V^{\pi_\theta}(s)$  and  $Q^{\pi_\theta}(s, a)$
- ▶ Using two function approximators and two parameter vectors:

$$V_v(s) \approx V^{\pi_\theta}(s)$$

$$Q_w(s, a) \approx Q^{\pi_\theta}(s, a)$$

$$A(s, a) = Q_w(s, a) - V_v(s)$$

- ▶ And updating both value functions by e.g. TD learning

# Dueling Networks

- ▶ Split Q-network into two channels
- ▶ Action-independent value function  $V(s, v)$
- ▶ Action-dependent advantage function  $A(s, a, w)$

$$Q(s, a) = V(s, v) + A(s, a, \mathbf{w})$$

- ▶ Advantage function is defined as:

$$A^{\pi}(s, a) = Q^{\pi}(s, a) - V^{\pi}(s).$$

# Estimating the Advantage Function

- ▶ For the true **value function**  $V^{\pi_\theta}(s)$  the TD error:

$$\delta^{\pi_\theta} = r + \gamma V^{\pi_\theta}(s') - V^{\pi_\theta}(s)$$

is an **unbiased estimate** of the advantage function:

$$\begin{aligned}\mathbb{E}_{\pi_\theta} [\delta^{\pi_\theta} | s, a] &= \mathbb{E}_{\pi_\theta} [r + \gamma V^{\pi_\theta}(s') | s, a] - V^{\pi_\theta}(s) \\ &= Q^{\pi_\theta}(s, a) - V^{\pi_\theta}(s) \\ &= A^{\pi_\theta}(s, a)\end{aligned}$$

- ▶ So we can use the TD error to compute the **policy gradient**

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) \delta^{\pi_\theta}]$$

- ▶ Remember the policy gradient

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) A^{\pi_\theta}(s, a)]$$

# Estimating the Advantage Function

- ▶ For the true **value function**  $V^{\pi_\theta}(s)$  the TD error:

$$\delta^{\pi_\theta} = r + \gamma V^{\pi_\theta}(s') - V^{\pi_\theta}(s)$$

is an **unbiased estimate** of the advantage function

- ▶ So we can use the TD error to compute the **policy gradient**

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) \delta^{\pi_\theta}]$$

- ▶ In practice we can use **an approximate TD error**

$$\delta_v = r + \gamma V_v(s') - V_v(s)$$

- ▶ This approach only requires one set **of critic parameters**  $v$

# Critic

- ▶ Critic can estimate **value function**  $V_v(s)$  from various targets. For example, from previous lectures:

- ▶ For MC, the target is the return  $G_t$

$$\Delta v = \alpha(G_t - V_v(s_t)) \nabla_v V_v(s)$$

- ▶  $V_v$  can be a **deep neural network** with parameters  $v$ .

- ▶ For TD(0), the target is the TD target  $r + \gamma V(s')$

$$\Delta v = \alpha(r_{t+1} + \gamma V_v(s_{t+1}) - V_v(s_t)) \nabla_v V_v(s)$$

- ▶ So critic is updated to minimize MSE w.r.t. target, given by MC or TD(0)

$$(V^{\pi_\theta}(s_t) - V_v(s_t))^2$$

# Actor

- ▶ The **policy gradient** can also be estimated as follows:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \textcolor{red}{A}^{\pi_{\theta}}(s, a)]$$

- ▶ Monte-Carlo policy gradient uses error from complete return

$$\Delta\theta = \alpha(G_t - V_v(s_t)) \nabla_{\theta} \log \pi_{\theta}(s_t, a_t)$$

- ▶ **Actor-critic policy gradient** uses the one-step TD error

$$\Delta\theta = \alpha(\textcolor{red}{r} + \gamma \textcolor{red}{V}_v(\textcolor{red}{s}_{t+1}) - V_v(s_t)) \nabla_{\theta} \log \pi_{\theta}(s_t, a_t)$$

# Advantage Actor-Critic Algorithm

## One-step Actor-Critic (episodic)

Input: a differentiable policy parameterization  $\pi(a|s, \boldsymbol{\theta}), \forall a \in \mathcal{A}, s \in \mathcal{S}, \boldsymbol{\theta} \in \mathbb{R}^n$

Input: a differentiable state-value parameterization  $\hat{v}(s, \mathbf{w}), \forall s \in \mathcal{S}, \mathbf{w} \in \mathbb{R}^m$

Parameters: step sizes  $\alpha > 0, \beta > 0$

Initialize policy weights  $\boldsymbol{\theta}$  and state-value weights  $\mathbf{w}$

Repeat forever:

    Initialize  $S$  (first state of episode)

$I \leftarrow 1$

    While  $S$  is not terminal:

$A \sim \pi(\cdot|S, \boldsymbol{\theta})$

        Take action  $A$ , observe  $S', R$

$\delta \leftarrow R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})$       (if  $S'$  is terminal, then  $\hat{v}(S', \mathbf{w}) \doteq 0$ )

$\mathbf{w} \leftarrow \mathbf{w} + \beta \delta \nabla_{\mathbf{w}} \hat{v}(S, \mathbf{w})$

$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha I \delta \nabla_{\boldsymbol{\theta}} \log \pi(A|S, \boldsymbol{\theta})$

$I \leftarrow \gamma I$

$S \leftarrow S'$

# So Far: Summary of PG Algorithms

- ▶ The policy gradient has many equivalent forms

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) G_t] && \text{REINFORCE} \\ &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q^w(s, a)] && \text{Q Actor-Critic} \\ &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) A^w(s, a)] && \text{Advantage Actor-Critic} \\ &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \delta] && \text{TD Actor-Critic}\end{aligned}$$

- ▶ Each leads a stochastic gradient ascent algorithm
- ▶ Critic uses **policy evaluation** (e.g. MC or TD learning) to estimate  $Q^{\pi}(s, a)$ ,  $A^{\pi}(s, a)$  or  $V^{\pi}(s)$



# Bias in Actor-Critic Algorithms

- ▶ Approximating the policy gradient introduces **bias**

$$Q_w(s, a) \approx Q^{\pi_\theta}(s, a)$$

- ▶ A biased policy gradient may not find the right solution
- ▶ Luckily, if we choose value function approximation carefully
- ▶ Then we can avoid introducing any bias
- ▶ i.e. we can still follow the exact policy gradient

# Compatible Function Approximation

- ▶ If the following two conditions are satisfied:

1. Value function approximator is compatible to the policy

$$\nabla_w Q_w(s, a) = \nabla_\theta \log \pi_\theta(s, a)$$

2. Value function parameters  $w$  minimize the mean-squared error

$$\varepsilon = \mathbb{E}_{\pi_\theta} [(Q^{\pi_\theta}(s, a) - Q_w(s, a))^2]$$

- ▶ Then the policy gradient is exact,

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)]$$

- ▶ Remember:

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) Q^{\pi_\theta}(s, a)]$$

# Proof

- ▶ If  $w$  is chosen to minimize mean-squared error, gradient of  $\varepsilon$  w.r.t.  $w$  must be zero,

$$\nabla_w \varepsilon = 0$$

$$\mathbb{E}_{\pi_\theta} [(Q^\theta(s, a) - Q_w(s, a)) \nabla_w Q_w(s, a)] = 0$$

$$\mathbb{E}_{\pi_\theta} [(Q^\theta(s, a) - Q_w(s, a)) \nabla_\theta \log \pi_\theta(s, a)] = 0$$

$$\mathbb{E}_{\pi_\theta} [Q^\theta(s, a) \nabla_\theta \log \pi_\theta(s, a)] = \mathbb{E}_{\pi_\theta} [Q_w(s, a) \nabla_\theta \log \pi_\theta(s, a)]$$

- ▶ So  $Q_w(s, a)$  can be substituted directly into the policy gradient,

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)]$$

# Alternative Policy Gradient Directions

- ▶ Gradient ascent algorithms can follow any **ascent direction**
- ▶ A good ascent direction can significantly speed convergence
- ▶ Also, a policy can often be **reparametrized** without changing action probabilities
- ▶ For example, increasing score of all actions in a softmax policy
- ▶ The vanilla gradient is sensitive to these reparametrizations

# Natural Policy Gradient

- ▶ The **natural policy gradient** is parametrization independent
- ▶ it finds ascent direction that is closest to vanilla gradient, when changing policy by a small, fixed amount

$$\nabla_{\theta}^{nat} \pi_{\theta}(s, a) = G_{\theta}^{-1} \nabla_{\theta} \pi_{\theta}(s, a)$$

- ▶ where  $G_{\theta}$  is the **Fisher information** matrix

$$G_{\theta} = \mathbb{E}_{\pi_{\theta}} \left[ \nabla_{\theta} \log \pi_{\theta}(s, a) \nabla_{\theta} \log \pi_{\theta}(s, a)^T \right]$$

# Natural Actor-Critic

- ▶ Using compatible function approximation,

$$\nabla_w A_w(s, a) = \nabla_\theta \log \pi_\theta(s, a)$$

- ▶ The **natural policy gradient** simplifies,

$$\begin{aligned}\nabla_\theta J(\theta) &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) A^{\pi_\theta}(s, a)] \\ &= \mathbb{E}_{\pi_\theta} \left[ \nabla_\theta \log \pi_\theta(s, a) \nabla_\theta \log \pi_\theta(s, a)^\top w \right] \\ &= G_\theta w\end{aligned}$$

$$\nabla_\theta^{\text{nat}} J(\theta) = w$$

Under linear model:

$$A^{\pi_\theta}(s, a) = \phi(s)^\top w$$

- ▶ i.e. update actor parameters in direction of critic parameters

# Summary of Policy Gradient Algorithms

- ▶ The policy gradient has many equivalent forms

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) G_t] \quad \text{REINFORCE}$$

$$= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q^w(s, a)] \quad \text{Q Actor-Critic}$$

$$= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) A^w(s, a)] \quad \text{Advantage Actor-Critic}$$

$$= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \delta] \quad \text{TD Actor-Critic}$$

$$= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \delta e] \quad \text{TD}(\lambda) \text{ Actor-Critic}$$

$$G_{\theta}^{-1} \nabla_{\theta} J(\theta) = w \quad \text{Natural Actor-Critic}$$

- ▶ Each leads a stochastic gradient ascent algorithm
- ▶ Critic uses **policy evaluation** (e.g. MC or TD learning) to estimate  $Q^{\pi}(s, a)$ ,  $A^{\pi}(s, a)$  or  $V^{\pi}(s)$

# Caption Generation with Visual Attention



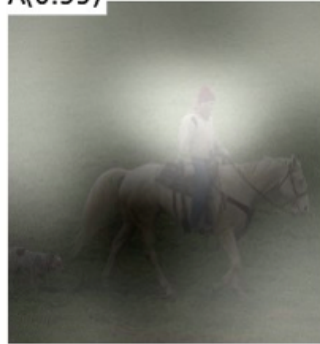
A man riding a horse  
in a field.



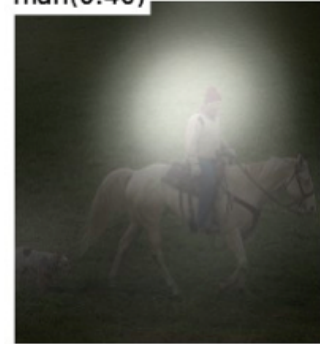
# Caption Generation with Visual Attention



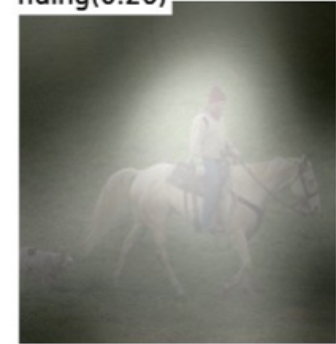
A(0.99)



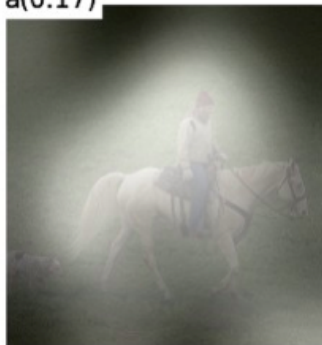
man(0.40)



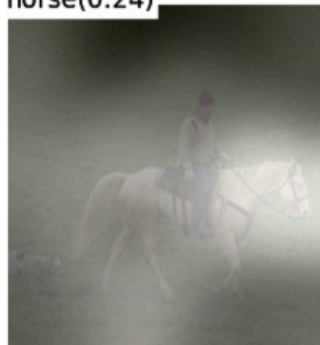
riding(0.26)



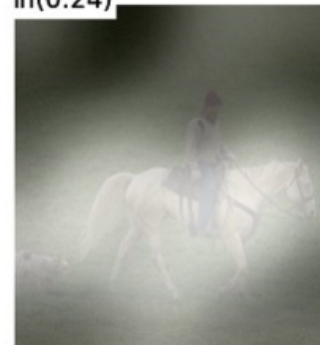
a(0.17)



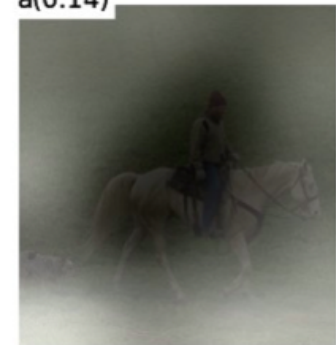
horse(0.24)



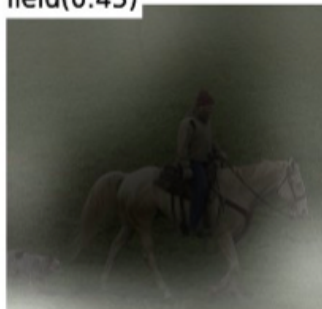
in(0.24)



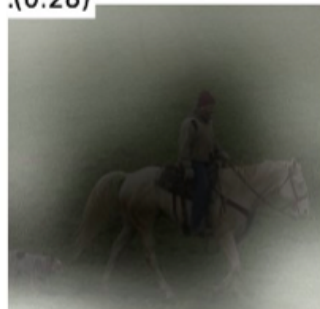
a(0.14)



field(0.43)



.(0.28)



A man riding a horse  
in a field.

# Caption Generation with Visual Attention



A woman is throwing a frisbee in a park.



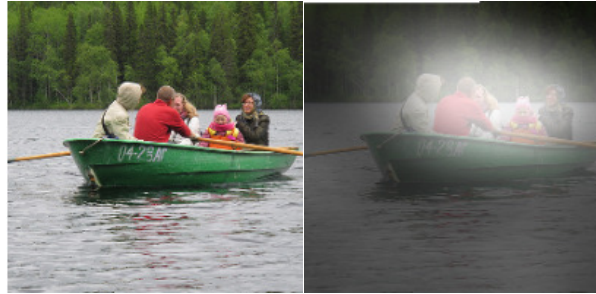
A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

# Improving Action Recognition

- Consider performing action recognition in a video:



- Instead of processing each frame, we can process only a small piece of each frame.

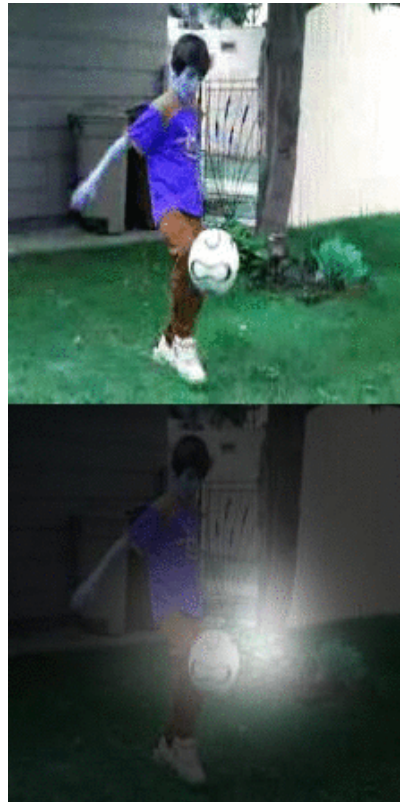


# Improving Action Recognition

Cycling



Soccer juggling



Horse back riding



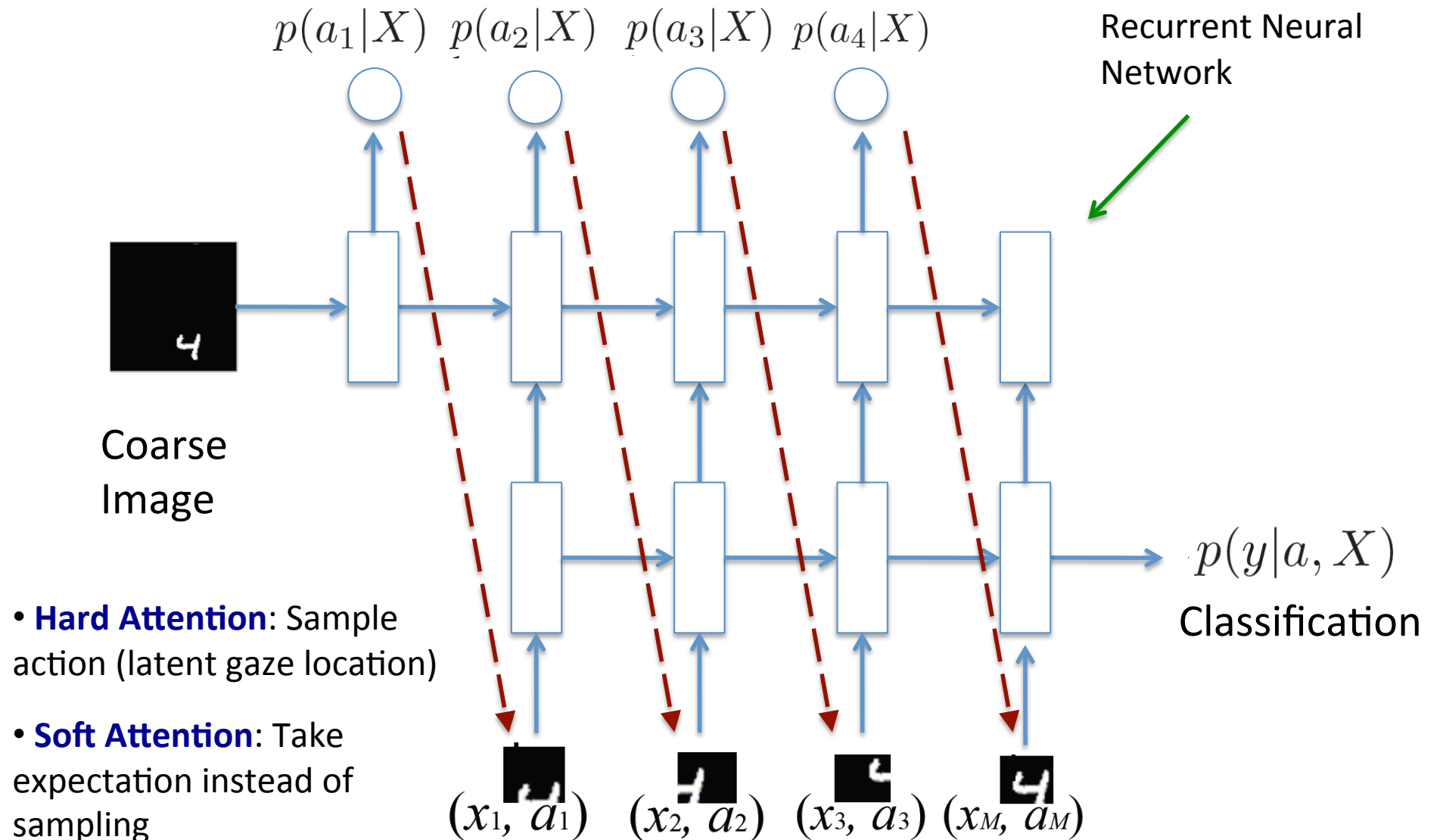
Basketball Shooting



# Recurrent Attention Model

Sample action:

$$\tilde{a}_1 \sim p(a_1|X) \quad \tilde{a}_2 \sim p(a_2|X)$$



# Model Setup

- We assume that we have a dataset with labels  $y$  for the supervised prediction task (e.g. object category).
- **Goal:** Learn an attention policy: The best locations to attend to are the ones which lead the model to predict the correct class.



# Model Definition

- We aim to maximize the probability of correct class by marginalizing over the actions (or latent gaze locations):

$$\mathcal{L} = \log p(y|X, W) = \log \sum_a p(a|X, W) p(y|a, X, W).$$

where

- $W$  is the set of parameters of the recurrent network.
- $a$  is a set of actions (latent gaze locations, scale).
- $X$ : is the input (e.g. image, video frame).

For clarity of presentation, I will sometimes omit conditioning on  $W$  or  $X$ . It should be obvious from the context.

# Variational Learning

- Previous approaches used variational lower bound:

$$\mathcal{LL} = \log \sum_a p(a|X, W) p(y|a, X, W) \geq \sum_a q(a|y, X) \log p(y, a|X, W) + \mathcal{H}[q] = \mathcal{F}.$$

- Here  $q(a|y, X)$  is some approximation to posterior over the gaze locations.
- In the case where  $q$  is the prior,  $q(a|y, X) = p(a|X, W)$ , the variational bound becomes:

$$\mathcal{F} = \sum_a p(a|X, W) \log p(y|a, X, W).$$

Ba et.al., ICLR 2015  
Mnih et.al., NIPS 2014



# REINFORCE

$$\mathcal{F} = \sum_a p(a|X, W) \log p(y|a, X, W).$$

- Derivatives w.r.t model parameters:

$$\frac{\partial \mathcal{F}}{\partial W} = \sum_a p(a|X, W) \left[ \frac{\partial \log p(y|a, X, W)}{\partial W} + \underbrace{\log p(y|a, X, W)}_{\text{Very bad term as it is unbounded.}} \frac{\partial \log p(a|X, W)}{\partial W} \right].$$

Very bad term as it is unbounded.  
Introduces high variance in the estimator.

- Need to introduce heuristics (e.g. replacing this term with a 0/1 discrete indicator function, which leads to REINFORCE algorithm of Williams, 1992).

Ba et.al., ICLR 2015  
Mnih et.al., NIPS 2014

# REINFORCE

$$\mathcal{F} = \sum_a p(a|X, W) \log p(y|a, X, W).$$

- Derivatives w.r.t model parameters:

$$\frac{\partial \mathcal{F}}{\partial W} = \sum_a p(a|X, W) \left[ \frac{\partial \log p(y|a, X, W)}{\partial W} + \log p(y|a, X, W) \frac{\partial \log p(a|X, W)}{\partial W} \right].$$

- The stochastic estimator of the gradient is given by:

$$\frac{\partial \mathcal{F}}{\partial W} \approx \frac{1}{M} \sum_{m=1}^M \left[ \frac{\partial \log p(y|\tilde{a}^m, X, W)}{\partial W} + \log p(y|\tilde{a}^m, X, W) \frac{\partial \log p(\tilde{a}^m|X, W)}{\partial W} \right].$$

where we draw M actions from the prior:  $\tilde{a}^m \sim p(a|X, W)$ .

# MNIST Attention Demo

- Actions contain:
  - Location: 2-d Gaussian latent variable
  - Scale: 3-way softmax over 3 different scales

