# ACTION-VALUE APPROXIMATION - ON POLICY

- CANNOT INTO FULL, TABULAR VALUE FCNS FOR TASKS WITH LARGE STATE AND ACTION SPACES, OR CONTINUOUS SPACES.
- IDEA! GENERALIZE FROM KNOWN, AVAILABLE VALUES ⟶ FUNCTION APPROXIMATION, REGRESSION, SUPERVISED LEARNING **ANY SUPERVISED LEARNING IS ✓**
- CAN USE ANNS OR OTHER TECHS ⟼ BACKUPS AREN'T TRIVIAL SHIFTS ANYMORE BUT ARBITRARILY COMPLEX
- IDEA: USE $S \longmapsto V$ BACKUP PAIRS AS SUPERVISED INSTANCES. OUTPUT IS ESTIMATED VALUE FUNCTION
- SUITABLE SV ARE THOSE CONVENIENTLY ALLOWING ONLINE LEARNING OF NONSTATIONARY STUFF.
- PERFORMANCE METRIC: RMSE

$$RMSE(w) = \sqrt{\sum d(s)\left[v_\pi(s) - \hat{v}(s,w)\right]^2}$$

$d$ { IS DISTRIBUTION SPECIFYING RELATIVE IMPORTANCE OF ERROR IN DIFFERENT STATES
DISTRIBUTION WE DRAW TRAINING EXAMPLES FROM
DISTRIBUTION WE DO BACKUPS OF (OF STATES)
$d = \pi$ ⟶ ON-POLICY DISTRIBUTION, THAT OF FREQ. OF STATES ENCOUNTERED }

- UNCLEAR IF MINIMIZING RMSE COINCIDES WITH ULTIMATE GOAL OF FINDING BETTER POLICIES

## GRADIENT DESCENT

- $W$ WEIGHTS, $\hat{v}(s,w)$, $w_T$ WEIGHTS AT STEP $t$. ASSUME A NEW EXAMPLE PER STEP $(s, v_\pi(s))$ TRUE VALUE UNDER $\pi$

$$w_{T+1} = w_T + \alpha\left[v_\pi(s_T) - \hat{v}_\pi(s_T, w_T)\right]\nabla\hat{v}(s_T, w_T)$$

- **NO TRUE VALUES** BECAUSE NOISE OR OTHER SHIT. USE ESTIMATE / BACKUP $v_T$ INSTEAD, SAME FORM. STILL CONVERGES FOR DECREASING $\alpha$s
- $\underline{v_T = G_T}$ MC RETURN ⟶ CONVERGENCE OK   • $v_T = G_T^\lambda$ $\lambda$-RETURN IS NOT UNBIASED ESTIMATE, NO CONVERGENCE, STILL EFFECTIVE AND RELEVANT

- { $w_{T+1} = w_T + \alpha\delta_T e_T$
  $\delta_T = R_{T+1} + \gamma\hat{v}(s_T, w_T) - \hat{v}(s_T, w_T)$ ⟵ **BWD TO (λ)** GRADIENT WITH $e$ ELIGIBILITY TRACES   • USUALLY ANN + BACKPROP, OR LINEAR METHODS
  $e_T = \lambda\gamma e_{T-1} + \nabla\hat{v}(s_T, w_T)$ }

**LINEAR METHODS:** { $\hat{v}(s,w) = w^T x$   • $x(s)$ FEATURE VECTOR, STATE IS IDENTIFIED WITH FEATURES !!! **LINEAR REGRESSION!**
$\nabla\hat{v}(s,w) = x(s)$
   • CONVERGENCE IS GUARANTEED AND BUT TO $w_\infty$, NOT ACTUAL $w^*$  ↳ MIN POSSIBLE ERROR
   • FEATURE SELECTION IS KEY }

**ON-POLICY + LINEAR!**
$$RMSE(w_\infty) \le \frac{1-\gamma\lambda}{1-\gamma}RMSE(w^*)$$

- COARSE CODING: FEATURES WITH 'RECEPTIVE FIELDS' THAT OVERLAP. (LINCOMB?) PRECISION RESOLUTION/GENERALIZATION TRADEOFF ON SIZE OF RECEPTIVE FIELDS AND NO OF FEATURES. ⟶ **BANDWIDTH**

- TILE CODING : COARSE, WITH ONLY BINARY FEATURES WITH NON-OVERLAPPING R.F. TILINGS NOT NECESSARILY UNIFORM GRIDS. COMPUTATION BECOMES VERY EFFICIENT BECAUSE WE CAN COUNT/SUM INSTEAD THAN MULTIPLY. ALSO HASHING.

- RBF CODING : GAUSSIAN KERNEL. SMOOTHER, DIFFERENTIABLE APPROXIMATIONS ARE PRODUCED. CAN ALSO LEARN $\mu, \Sigma$ FOR GREAT JUSTICE.

- KANERVA CODING: DECOUPLE DIMENSIONALITY/COMPLEXITY OF STATE SPACE FROM THAT OF TARGET FUNCTION. FOR RESILIENCE WRT CURSE OF DIMENSIONALITY ⟶ USE PROTOTYPES (KERNELS?) AND EXPRESS STATES WRT CLOSENESS TO PROTOTYPES **EXAMPLE:** BINARY SPACE, HAMMING DISTANCE

## CONTROL WITH FUNCTION APPROXIMATION

- USUAL GPI APPROACH
- ACTION-VALUE PREDICTION: $\hat{q} \approx q_\pi(s_T, A_T) \longrightarrow Q_T$: MC RETURN, SARSA RETURN, ETC... $w_{T+1} = w_T + \alpha\left[q_T - \hat{q}(s_T, A_T, w_T)\right]\nabla\hat{q}(s_T, A_T, w_T)$
- POLICY IMPROVEMENT, **ACTION SELECTION:**
  - DISCRETE ACTION SET, NOT TOO LARGE ⟶ USUAL TECHNIQUES OK : FOR EACH $a$, COMPUTE $\hat{q}$, AND FIND GREEDY ACTION
  - CONTINUOUS OR LARGE DISCRETE ACTION SPACES ⟶ NO CLEAR SOLUTION / ONGOING RESEARCH

  { **ON-POLICY** PI W/ SOFT APPROX, $\varepsilon$-GREEDY, BEHAVE WITH POLICY } SAME
  { **OFF-POLICY** PI W/ GREEDY POLICY, SELECTION W/ ARBITRARY POLICY }

- TRACES : CAN USE ANY VARIANT. NO TRACE FOR STATES, BUT FOR WEIGHTS ⟶ TREAT FEATURES AS STATES AND DO TRACES ON THEM. OPTIONAL CLEARING OF TRACES OF NONSELECTED ACTIONS

## BOOTSTRAPPING

- BOOTSTRAPPING : BETTER AT SOLVING PROBLEMS. NOT CLEAR WHY.
- NONBOOTSTRAPPING : BETTER AT MINIMIZING RMSE.