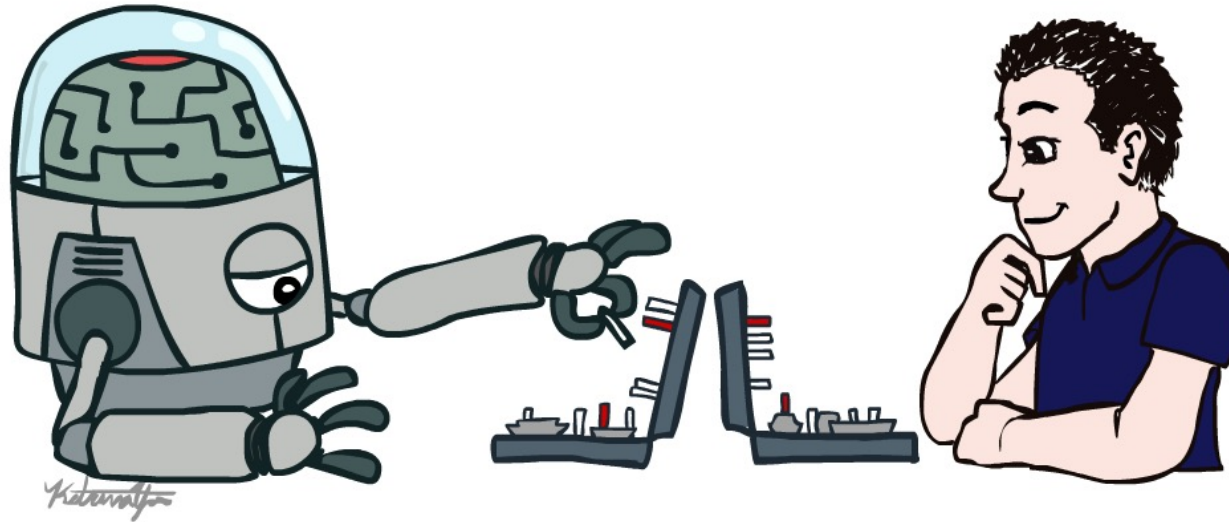# CS 188: Artificial Intelligence

## Final Exam Review

Instructor: Nicholas Tomlin

University of California, Berkeley

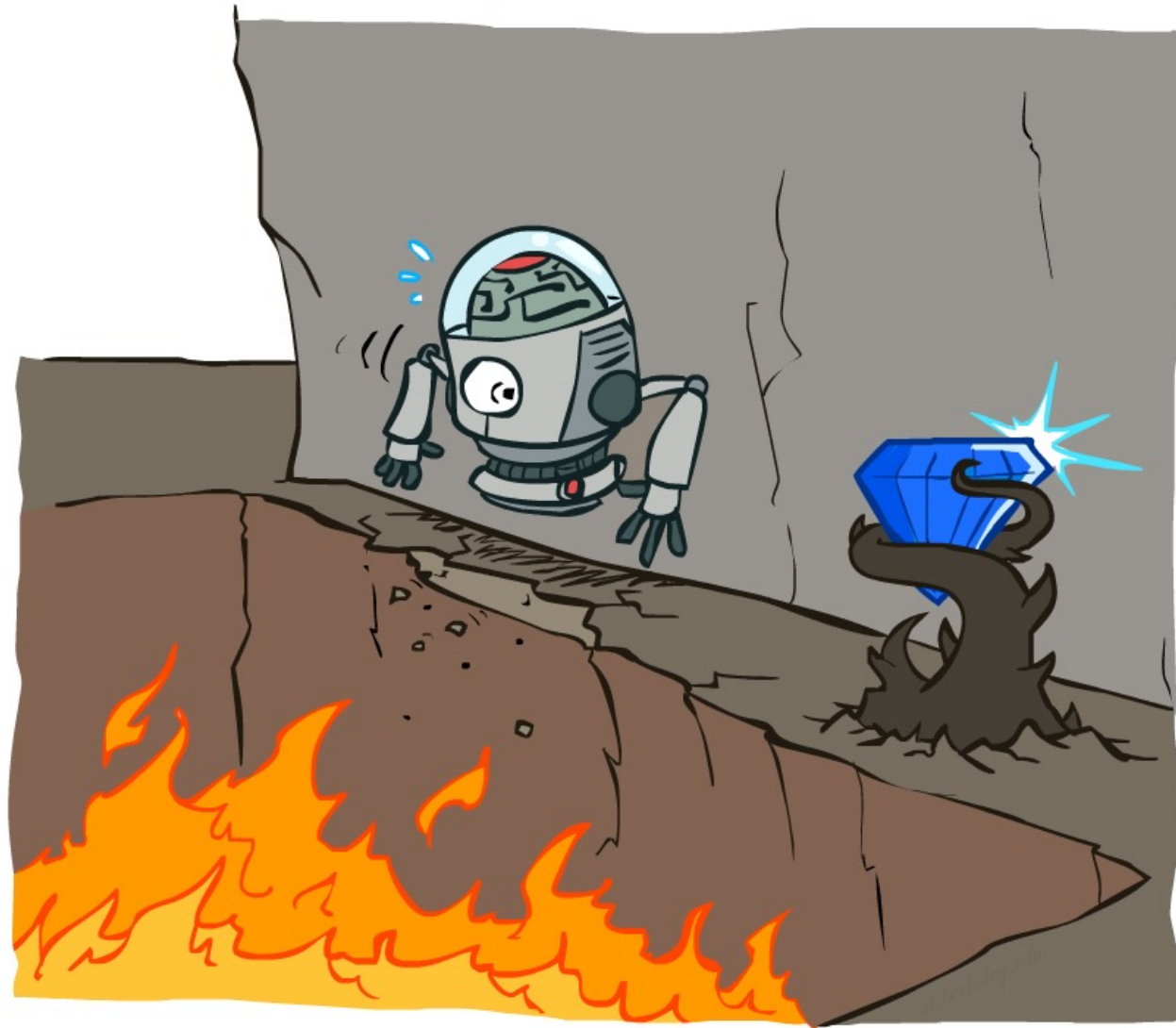(slides adapted from Dan Klein, Pieter Abbeel, Anca Dragan, Stuart Russell)

# Announcements

- All assignments **must** be submitted by 11:59PM tonight; no late submissions, no grace period beyond this point
- Exam tomorrow (Thursday, August 10th):
  - Wheeler 150
  - 7-10PM, but please show up no later than 6:45PM
  - Get enough sleep, drink enough water, etc.
- Get +1% extra credit on the exam by filling out course evaluations: https://course-evaluations.berkeley.edu

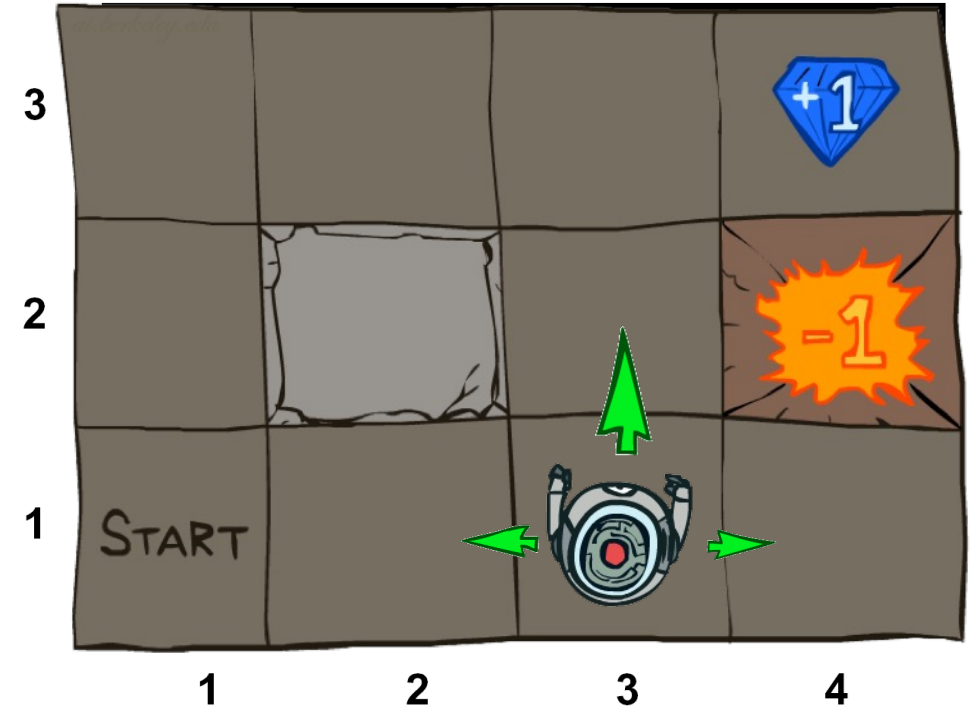- Today's plan: cover as much material as possible, focused on the second half of the course.

# Markov Decision Processes

# Markov Decision Processes

o An MDP is defined by:
  o A set of states s ∈ S
  o A set of actions a ∈ A
  o A transition function T(s, a, s')
    o Probability that a from s leads to s', i.e., P(s' | s, a)
    o Also called the model or the dynamics
  o A reward function R(s, a, s')
    o Sometimes just R(s) or R(s')
  o A start state
  o Maybe a terminal state

o We care about:
  o Policy = choice of actions for each state
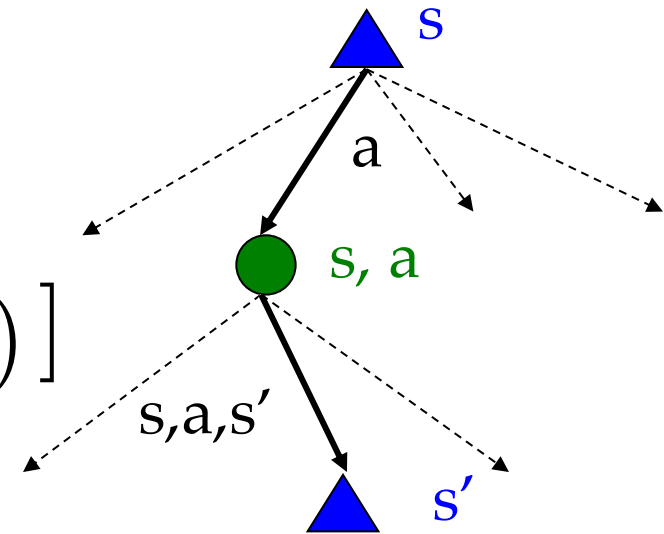  o Utility = sum of (discounted) rewards

# Values of States: Bellman Equation

o Recursive definition of value:

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

$$Q^*(s,a) = \sum_{s'} T(s,a,s')[R(s,a,s') + \gamma V^*(s')]$$

$$V^*(s) = \max_a \sum_{s'} T(s,a,s')[R(s,a,s') + \gamma V^*(s')]$$

# Value Iteration

o Bellman equations characterize the optimal values:

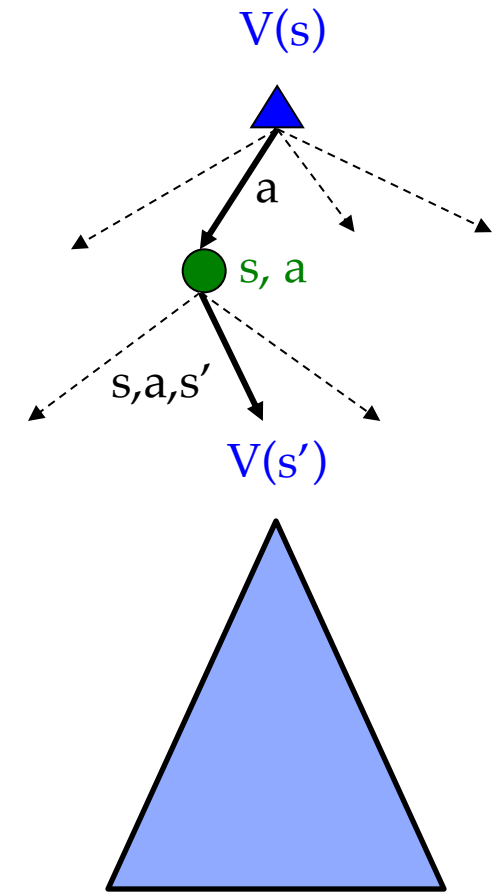$$V^*(s) = \max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V^*(s') \right]$$

o Value iteration computes them:

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V_k(s') \right]$$

"Bellman Update"

o Value iteration is just a fixed point solution method

V(s)

s, a

s,a,s'

V(s')

# Policy Extraction from Values

○ Let's imagine we have the optimal values V*(s)

○ How should we act?

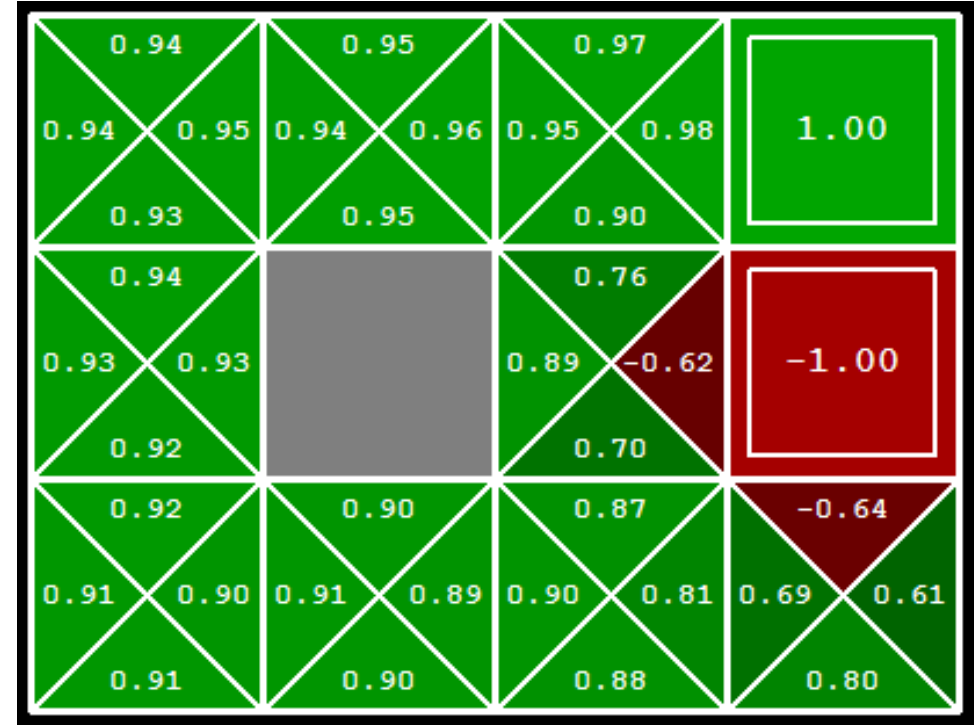　○ It's not obvious!

○ We need to do a mini-expectimax (one step)

$$\pi^*(s) = \arg\max_a \sum_{s'} T(s, a, s')[R(s, a, s') + \gamma V^*(s')]$$

○ This is called policy extraction, since it gets the policy implied by the values

# Policy Extraction from Q-Values

o Let's imagine we have the optimal q-values:

o How should we act?

    o Completely trivial to decide!

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

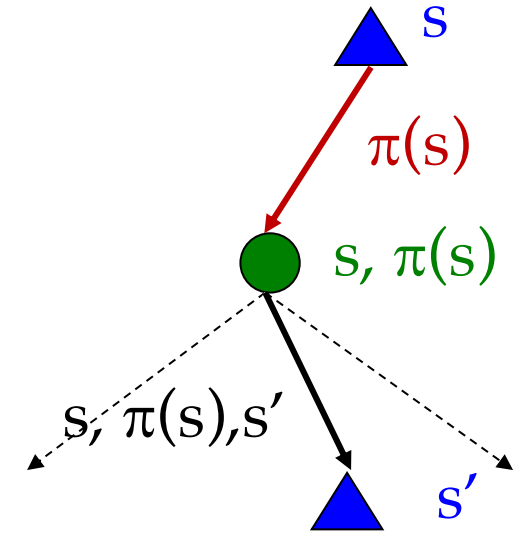o Important lesson: actions are easier to select from q-values than values!

# Policy Evaluation

- How do we calculate the V's for a fixed policy π?

- Idea 1: Turn recursive Bellman equations into updates (like value iteration)

$$V_0^\pi(s) = 0$$

$$V_{k+1}^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s')[R(s, \pi(s), s') + \gamma V_k^\pi(s')]$$

- Efficiency: $O(S^2)$ per iteration

- Idea 2: Without the maxes, the Bellman equations are just a linear system
  - Solve the system of equations

s

π(s)

s, π(s)

s, π(s),s'

s'

# Policy Iteration

o Evaluation: For fixed current policy π, find values with policy evaluation:

   o Iterate until values converge:

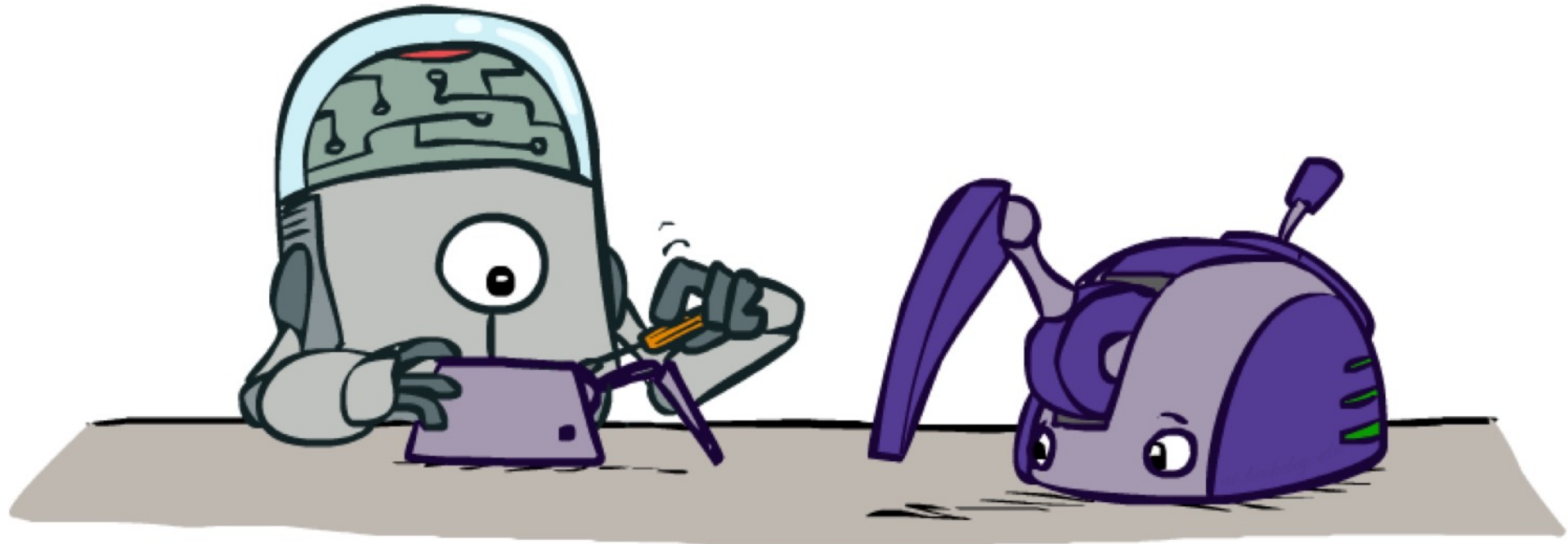$$V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s') \left[ R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s') \right]$$

o Improvement: For fixed values, get a better policy using policy extraction

   o One-step look-ahead:

$$\pi_{i+1}(s) = \arg\max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V^{\pi_i}(s') \right]$$

# Reinforcement Learning

# Map of Reinforcement Learning

## Known MDP: Offline Solution

| Goal | Technique |
|------|-----------|
| Compute V*, Q*, $\pi$* | Value / policy iteration |
| Evaluate a fixed policy $\pi$ | Policy evaluation |

## Unknown MDP: Model-Based

| Goal | Technique |
|------|-----------|
| Compute V*, Q*, $\pi$* | VI/PI on approx. MDP |
| Evaluate a fixed policy $\pi$ | PE on approx. MDP |

## Unknown MDP: Model-Free

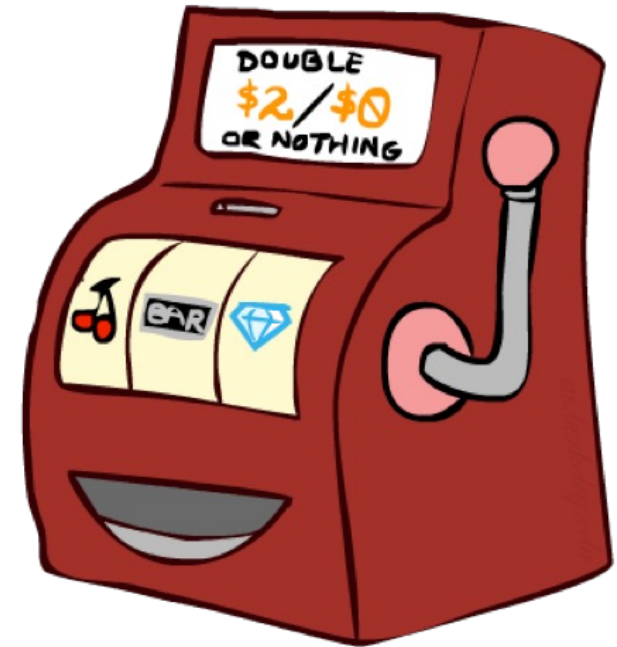| Goal | Technique |
|------|-----------|
| Compute V*, Q*, $\pi$* | Q-learning |
| Evaluate a fixed policy $\pi$ | Value Learning |

# Direct Evaluation
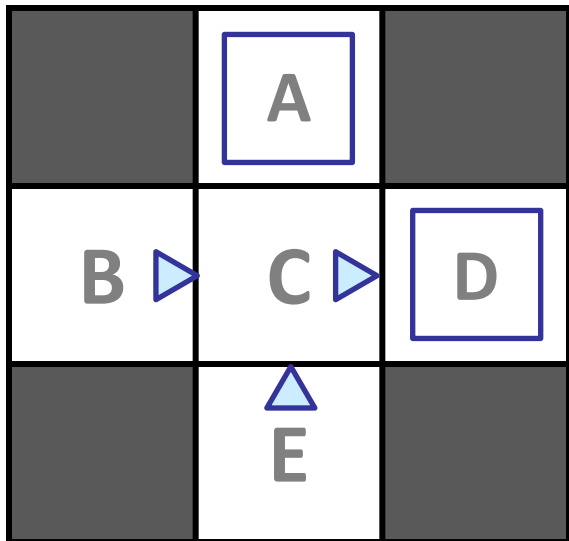
o Goal: Compute values for each state under π

o Idea: Average together observed sample values

 o Act according to π

 o Every time you visit a state, write down what the sum of discounted rewards turned out to be

 o Average those samples

o This is called direct evaluation

# Direct Evaluation

## Input Policy π



*Assume: γ = 1*

## Observed Episodes (Training)

### Episode 1

B, east, C, -1
C, east, D, -1
D, exit,  x, +10

### Episode 2

B, east, C, -1
C, east, D, -1
D, exit,  x, +10

### Episode 3
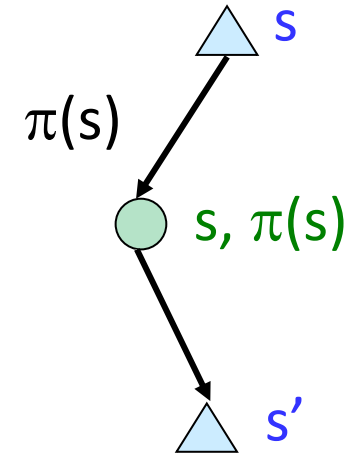
E, north, C, -1
C, east,   D, -1
D, exit,    x, +10

### Episode 4

E, north, C, -1
C, east,   A, -1
A, exit,    x, -10

## Output Values

# Temporal Difference Learning

○ Big idea: learn from every experience!

    ○ Update V(s) each time we experience a transition (s, a, s', r)

    ○ Likely outcomes s' will contribute updates more often

○ Temporal difference learning of values

    ○ Policy still fixed, still doing evaluation!

    ○ Move values toward value of whatever successor occurs: running average

$\pi(s)$

s

s, $\pi(s)$

s'

Sample of V(s):     $sample = R(s, \pi(s), s') + \gamma V^\pi(s')$

Update to V(s):     $V^\pi(s) \leftarrow (1 - \alpha)V^\pi(s) + (\alpha)sample$

Same update:     $V^\pi(s) \leftarrow V^\pi(s) + \alpha(sample - V^\pi(s))$

# Q-Learning

○ Q-Learning: sample-based Q-value iteration

$$Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma \max_{a'} Q_k(s', a') \right]$$

○ Learn Q(s,a) values as you go

  ○ Receive a sample (s,a,s',r)

  ○ Consider your old estimat $Q(s, a)$

  ○ Consider your new sample estimate:

  $$sample = R(s, a, s') + \gamma \max_{a'} Q(s', a')$$

  ○ Incorporate the new estimate into a running average:

  $$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) \, [sample]$$



Q-VALUES AFTER 1000 EPISODES

# Approximate Q-Learning

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \ldots + w_n f_n(s, a)$$

o Q-learning with linear Q-functions:
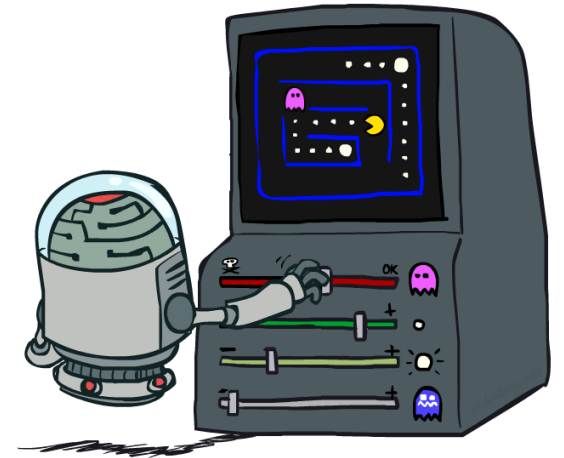
transition $= (s, a, r, s')$

difference $= \left[ r + \gamma \max_{a'} Q(s', a') \right] - Q(s, a)$

$Q(s, a) \leftarrow Q(s, a) + \alpha \, [\text{difference}]$      Exact Q's

$w_i \leftarrow w_i + \alpha \, [\text{difference}] \, f_i(s, a)$      Approximate Q's

o Intuitive interpretation:
  o Adjust weights of active features
  o E.g., if something unexpectedly bad happens, blame the features that were on: disprefer all states with that state's features

o Formal justification: online least squares

# How to Explore?

o Several schemes for forcing exploration
  o Simplest: random actions ($\varepsilon$-greedy)
    o Every time step, flip a coin
    o With (small) probability $\varepsilon$, act randomly
    o With (large) probability $1-\varepsilon$, act on current policy

  o Problems with random actions?
    o You do eventually explore the space, but keep thrashing around once learning is done
    o One solution: lower $\varepsilon$ over time
    o Another solution: exploration functions

# Exploration Functions

o **When to explore?**
  o Random actions: explore a fixed amount
  o Better idea: explore areas whose badness is no
    (yet) established, eventually stop exploring

o **Exploration function**
  o Takes a value estimate u and a visit count n, an
    returns an optimistic utility, e. $f(u,n) = u + k/n$

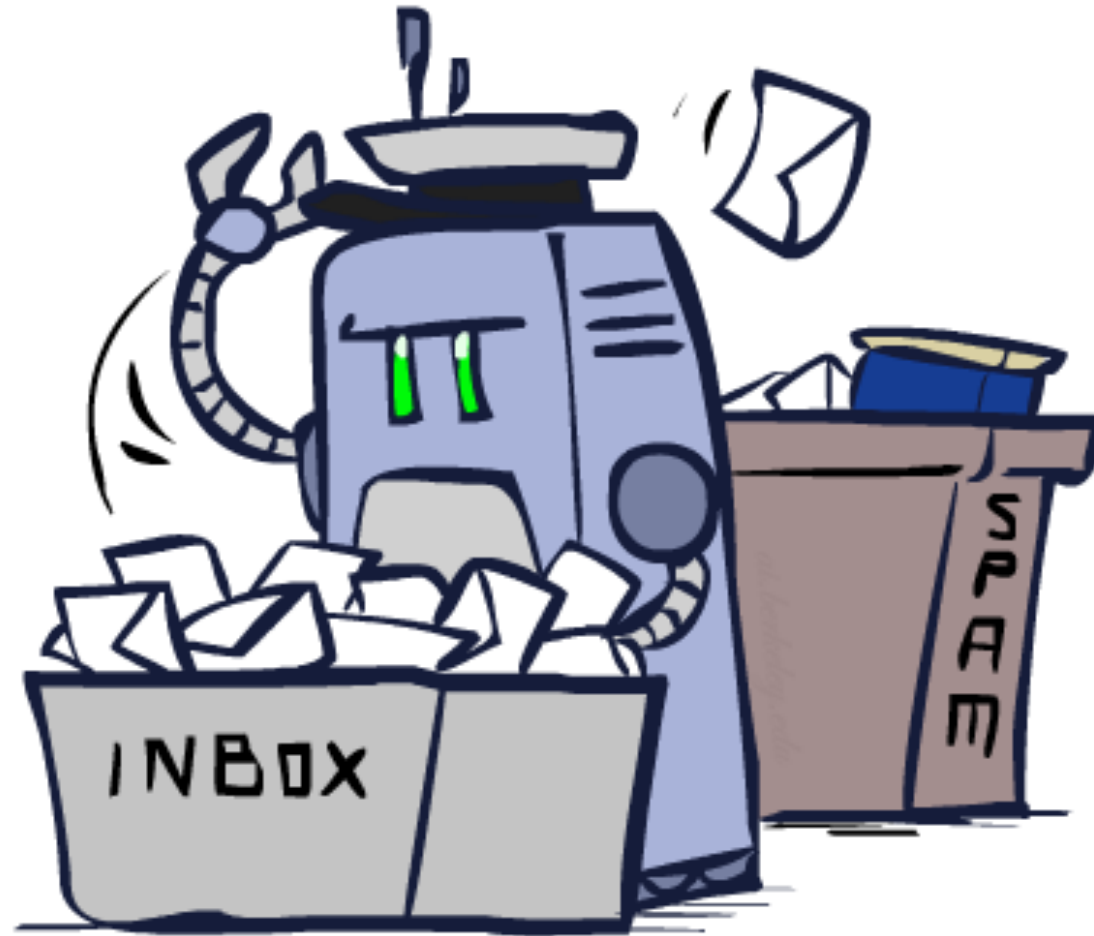  Regular Q-Update:    $Q(s,a) \leftarrow_\alpha R(s,a,s') + \gamma \max_{a'} Q(s',a')$

  Modified Q-Update:  $Q(s,a) \leftarrow_\alpha R(s,a,s') + \gamma \max_{a'} f(Q(s',a'), N(s',a'))$

  o Note: this propagates the "bonus" back to states that lead to unknown states
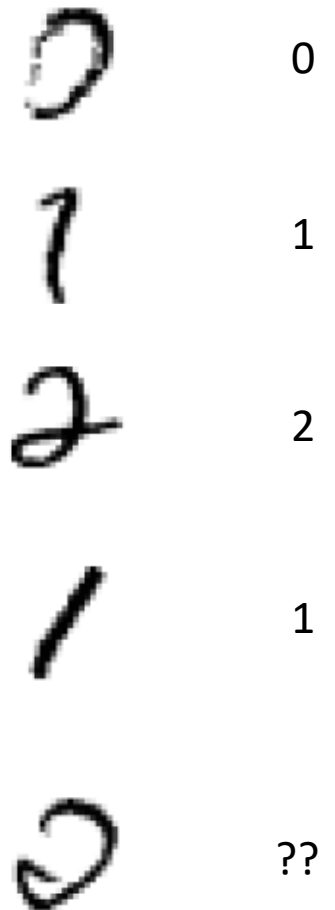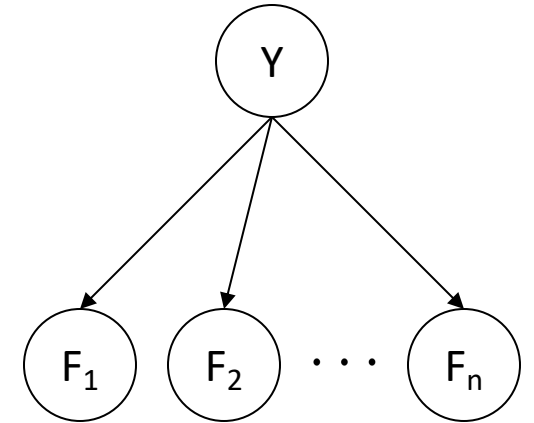    as well!

# Machine Learning

# Example: Digit Recognition

o Input: images / pixel grids

o Output: a digit 0-9

o Setup:
   o Get a large collection of example images, each labeled with a digit
   o Note: someone has to hand label all this data!
   o Want to learn to predict labels of new, future digit images

o Features: The attributes used to make the digit decision
   o Pixels: (6,8)=ON
   o Shape Patterns: NumComponents, AspectRatio, NumLoops
   o …
   o Features are increasingly induced rather than crafted

0

1

2

1

??

# Naïve Bayes for Digits

o Naïve Bayes: Assume all features are independent effects of the label

o Simple digit recognition version:



- o One feature (variable) $F_{ij}$ for each grid position $\langle i,j \rangle$
- o Feature values are on / off, based on whether intensity

  is more or less than 0.5 in underlying image
- o Each input maps to a feature vector, e.g.

$$\to \langle F_{0,0} = 0 \ \ F_{0,1} = 0 \ \ F_{0,2} = 1 \ \ F_{0,3} = 1 \ \ F_{0,4} = 0 \ \ \ldots F_{15,15} = 0 \rangle$$

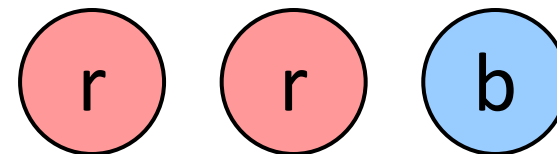- o Here: lots of features, each is binary valued

$$P_{\mathsf{ML}}(x) = \frac{\mathsf{count}(x)}{\mathsf{total\ samples}}$$

o Naïve Bayes model: $\quad P(Y|F_{0,0} \ldots F_{15,15}) \propto P(Y) \prod_{i,j} P(F_{i,j}|Y)$

o **Conditional probabilities $P(F_{i,j} \mid Y)$ just come from counts in the training data**

# Deriving MLEs

- **Model**:

| X | red | blue |
|---|-----|------|
| $P_\theta(x)$ | $\theta$ | $1 - \theta$ |

- **Data**: draw $N$ balls. $N_r$ come up red, $N_b$ come up blue
  - Dataset: $D = \{x_1, \dots, x_n\}$
  - Ball draws are independent and identically distributed (i.i.d.):

$$P(D \mid \theta) = \prod_i P(x_i \mid \theta) = \prod_i P_\theta(x_i) = \theta^{N_r} \cdot (1 - \theta)^{N_b}$$

- **Maximum likelihood estimation**: find $\theta$ that maximizes $P(D \mid \theta)$

$$\theta = \operatorname*{argmax}_\theta P(D \mid \theta) = \operatorname*{argmax}_\theta \log P(D \mid \theta)$$

  - Approach: take derivative and set to 0

# Deriving MLEs

○ **Maximum likelihood estimation**: find $\theta$ that maximizes $P(D \mid \theta)$

$$\theta = \operatorname*{argmax}_{\theta} P(D \mid \theta) = \operatorname*{argmax}_{\theta} \log P(D \mid \theta)$$

$$\frac{\partial}{\partial \theta} \log P(D \mid \theta) = \frac{\partial}{\partial \theta}[N_r \log(\theta) + N_b \log(1 - \theta)]$$

$$= N_r \frac{\partial}{\partial \theta} \log(\theta) + N_b \frac{\partial}{\partial \theta} \log(1 - \theta)$$

$$= N_r \frac{1}{\theta} - N_b \frac{1}{1-\theta}$$

$$= 0$$

Multiply by $\theta(1 - \theta)$: $\quad N_r(1 - \theta) - N_b\theta = 0$

$$N_r - \theta(N_r + N_b) = 0$$

$$\hat{\theta} = \frac{N_r}{N_r + N_b}$$

# Regularization: Smoothing

○ Laplace's estimate:

    ○ Pretend you saw every outcome once more than you actually did

$$P_{LAP}(x) = \frac{c(x) + 1}{\sum_x [c(x) + 1]}$$

$$= \frac{c(x) + 1}{N + |X|}$$

○ **This is no longer a maximum likelihood estimate**



$$P_{ML}(X) = \left\langle \frac{2}{3}, \frac{1}{3} \right\rangle$$

$$P_{LAP}(X) = \left\langle \frac{3}{5}, \frac{2}{5} \right\rangle$$

# Binary Perceptron

- Start with weights $= 0$
- For each training instance:
  - Classify with current weights

$$y = \begin{cases} +1 & \text{if } w \cdot f(x) \geq 0 \\ -1 & \text{if } w \cdot f(x) < 0 \end{cases}$$

  - If correct (i.e., y=y*), no change!
  - If wrong: adjust the weight vector by adding or subtracting the feature vector. Subtract if y* is -1.
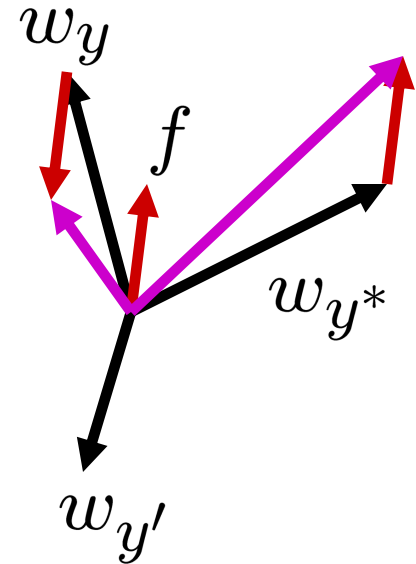
$$w = w + y^* \cdot f$$

# Multiclass Perceptron

- Start with all weights = 0
- Pick up training examples one by one
- Predict with current weights

$$y \; = \arg \max_y \; w_y \cdot f(x)$$

- If correct, no change!
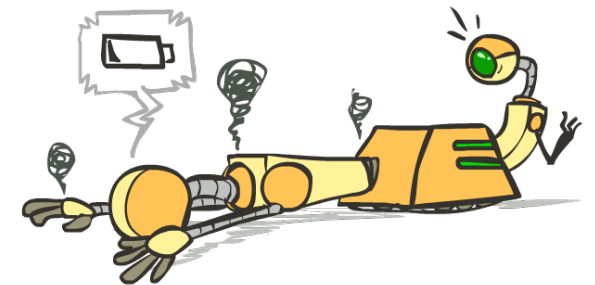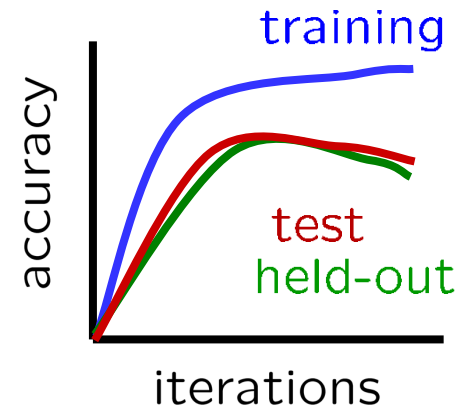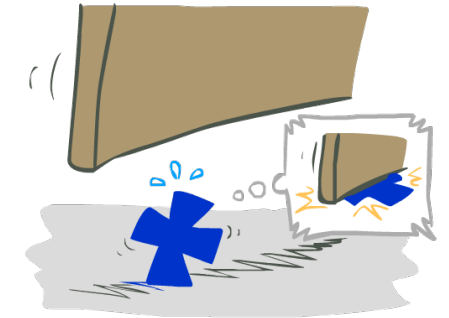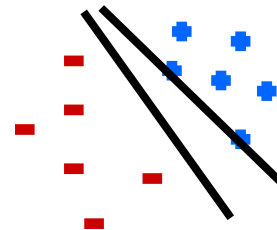- If wrong: lower score of wrong answer, raise score of right answer
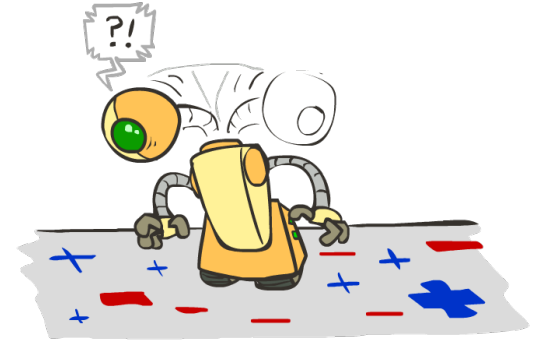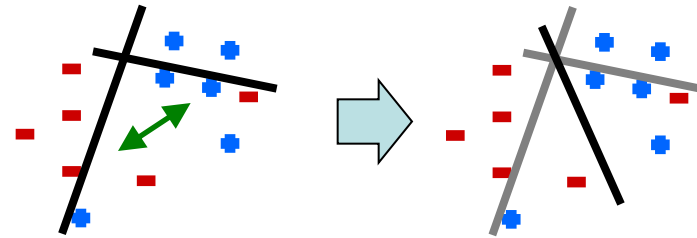
$$w_y = w_y - f(x)$$

$$w_{y*} = w_{y*} + f(x)$$

$w_y$

$f$

$w_{y*}$

$w_{y'}$

# Problems with the Perceptron

- Noise: if the data isn't separable, weights might thrash
  - Averaging weight vectors over time can help (averaged perceptron)

- Mediocre generalization: finds a "barely" separating solution

- Overtraining: test / held-out accuracy usually rises, then falls
  - Overtraining is a kind of overfitting

# Logistic Regression

○ Maximum likelihood estimation:

$$\max_w \quad ll(w) = \max_w \sum_i \log P(y^{(i)}|x^{(i)}; w)$$

with:
$$P(y^{(i)} = +1|x^{(i)}; w) = \frac{1}{1 + e^{-w \cdot f(x^{(i)})}}$$

$$P(y^{(i)} = -1|x^{(i)}; w) = 1 - \frac{1}{1 + e^{-w \cdot f(x^{(i)})}}$$
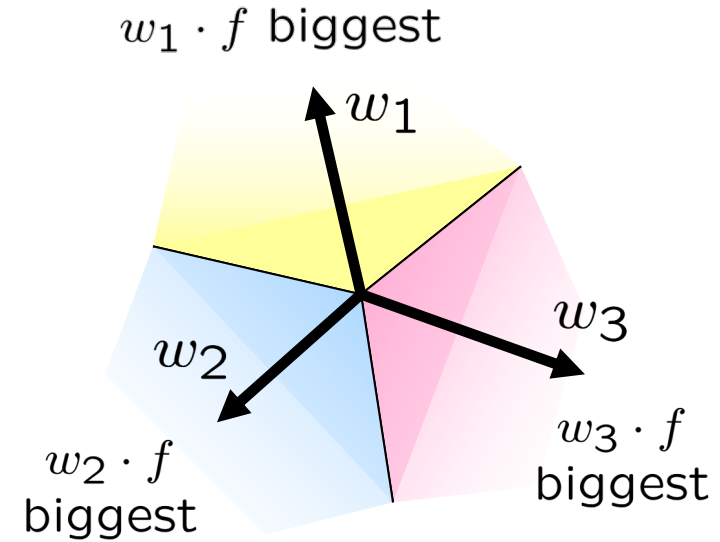
**Aside: linear regression ≠ logistic regression!**

# Multiclass Logistic Regression

- Recall Perceptron:
  - A weight vector for each class: $w_y$
  - Score (activation) of a class y: $w_y \cdot f(x)$
  - Prediction highest score wins $y = \arg\max_y \; w_y \cdot f(x)$

$w_1 \cdot f$ biggest

$w_1$

$w_2$

$w_3$

$w_2 \cdot f$
biggest

$w_3 \cdot f$
biggest

- How to make the scores into probabilities?

$$z_1, z_2, z_3 \rightarrow \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}, \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}, \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

original activations          softmax activations

# Batch Gradient Ascent

$$\max_{w} \quad ll(w) = \max_{w} \quad \underbrace{\sum_{i} \log P(y^{(i)}|x^{(i)}; w)}_{g(w)}$$

○ init $w$

○ for iter = 1, 2, …

$$w \leftarrow w + \alpha * \sum_{i} \nabla \log P(y^{(i)}|x^{(i)}; w)$$

# Stochastic Gradient Ascent

$$\max_{w} \; ll(w) = \max_{w} \; \sum_{i} \log P(y^{(i)}|x^{(i)}; w)$$

**Observation:** once gradient on one training example has been computed, might as well incorporate before computing next one

o init $w$

o for iter = 1, 2, …

  o pick random j

$$w \leftarrow w + \alpha * \nabla \log P(y^{(j)}|x^{(j)}; w)$$

# Mini-batch Gradient Ascent

$$\max_w \quad ll(w) = \max_w \quad \sum_i \log P(y^{(i)}|x^{(i)}; w)$$

**Observation:** gradient over small set of training examples (=mini-batch) can be computed in parallel, might as well do that instead of a single one

○ `init` $w$

○ `for iter = 1, 2, …`

  ○ `pick random subset of training examples J`

$$w \leftarrow w + \alpha * \sum_{j \in J} \nabla \log P(y^{(j)}|x^{(j)}; w)$$

# Beyond SGD: Second-Order Derivatives

## Newton's Method (in 1D):

- Want to optimize: $\max\limits_{\theta} f(\theta)$
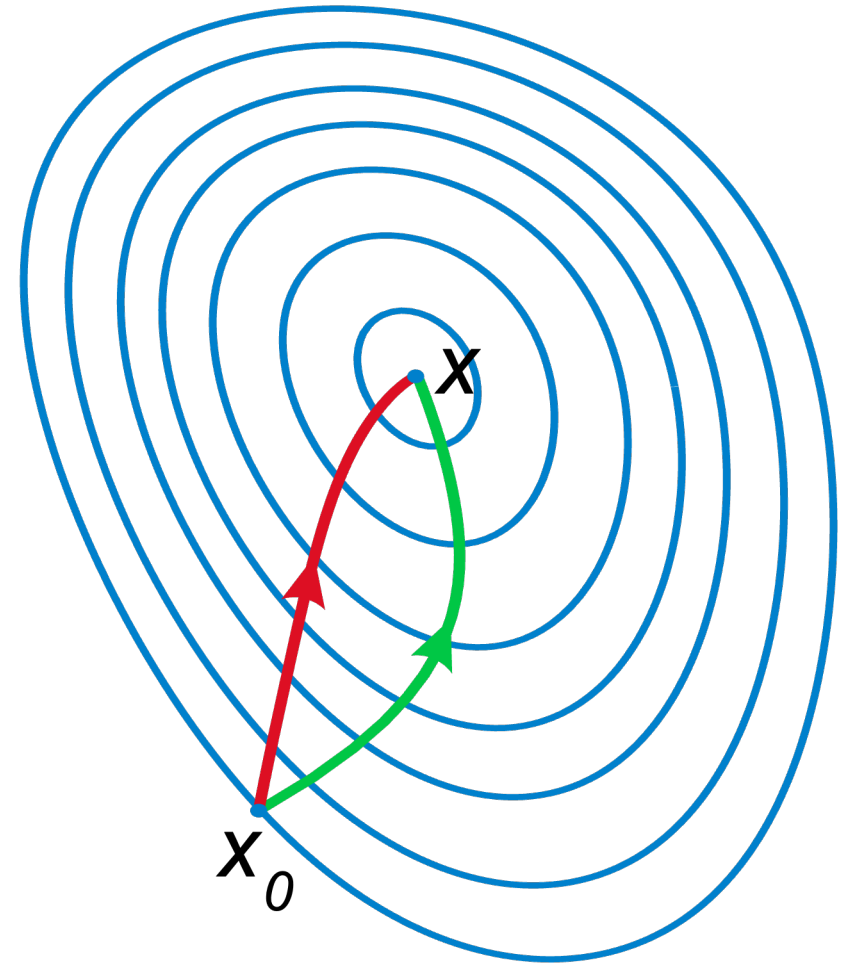
- Apply Taylor expansion:

$$f(\theta + h) = f(\theta) + f'(\theta)h + \frac{1}{2}f''(\theta)h^2$$

- Find value of $t$ that maximizes this:

$$0 = \frac{\partial}{\partial h}\left[f(\theta) + f'(\theta)h + \frac{1}{2}f''(\theta)h^2\right]$$
$$= f'(\theta) + f''(\theta)h$$

- Rearrange terms to get update:

$$h = -\frac{f'(\theta)}{f''(\theta)} \qquad \theta_{t+1} = \theta_t + h = \theta_t - \frac{f'(\theta)}{f''(\theta)}$$



**These update equations out of scope for final exam; but high-level concepts are in scope**

# Beyond SGD: Momentum

o Potential issues with vanilla SGD:

   o Can take a long time to converge if the learning rate is too low

   o Can bounce around in "ravines" without making much progress toward a local optimum
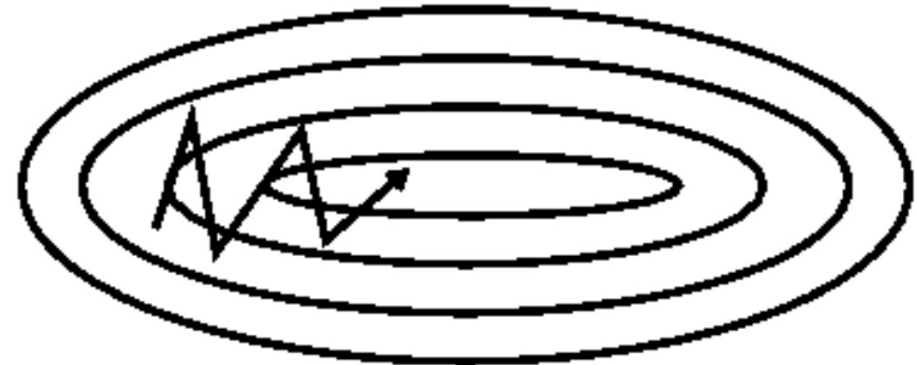


Image 2: SGD without momentum          Image 3: SGD with momentum

# Beyond SGD: Adaptive Learning Rates

o Recall: learning rates

- o Determines how much we update weights in the direction of the gradient
- o Often: want to set this in terms of how much it updates the weights
- o Often: want to lower learning rate over time (*learning rate scheduling*)

$$\theta_{t+1} = \theta_t - \boxed{\eta} \nabla_\theta f(\theta_t)$$
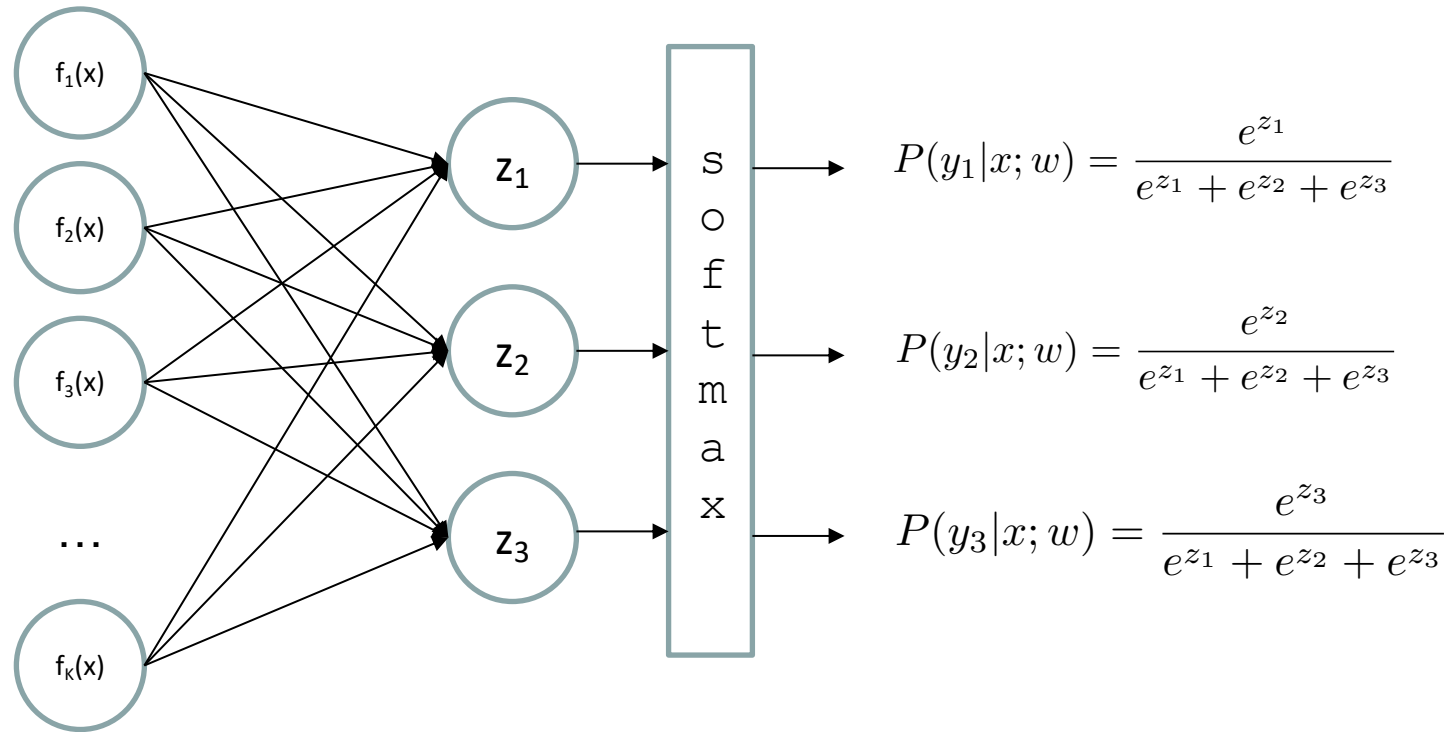
o Key idea: different learning rates for each parameter

- o We can make larger or smaller updates depending on how important a feature is
- o Small updates for frequent features; big updates for rare features
- o This idea underlies: Adagrad, RMSProp, Adam, etc.
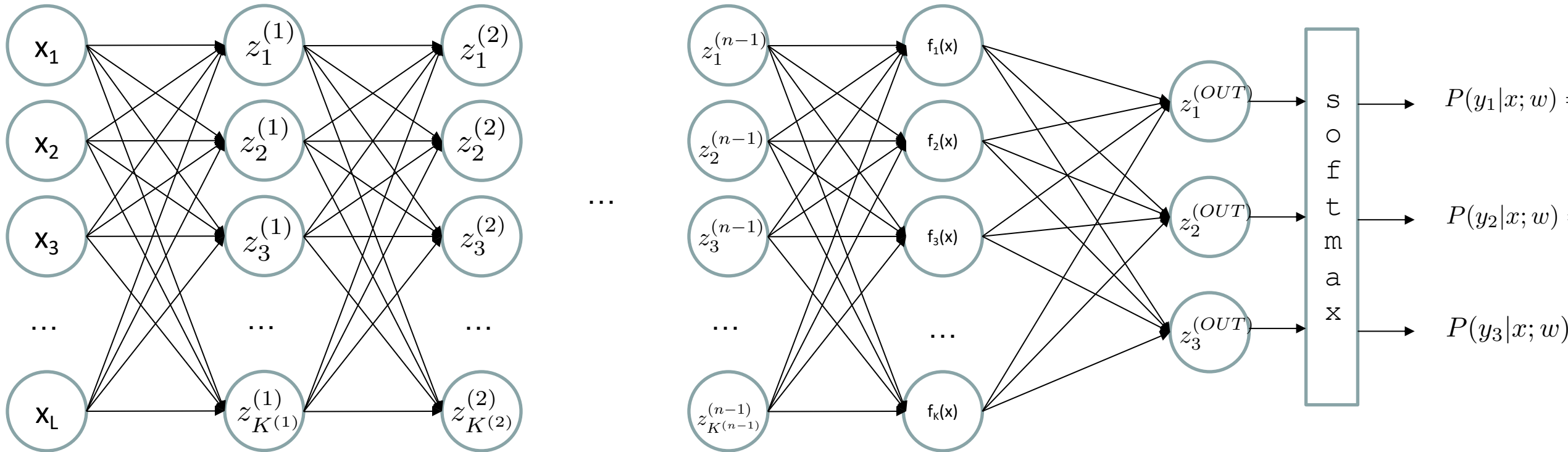
# Summary: Key Ideas in Optimization

o Gradient descent
  o Batch: update based on the whole dataset
  o SGD: update based on a single randomly chosen training example
  o Minibatch: update based on $k$ randomly chosen training examples

o More advanced approaches:
  o Second order optimization (e.g., Newton's method)
  o Momentum (Nesterov's accelerated gradient, Adam)
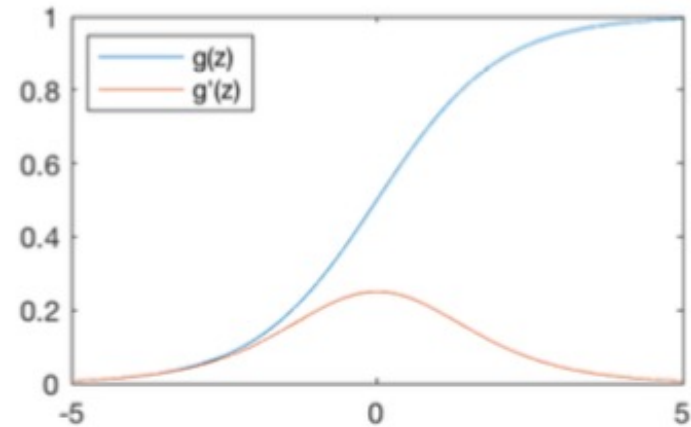  o Adaptive learning rates (Adagrad, RMSProp, Adam, etc.)

# Logistic Regression



$$P(y_1|x;w) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_2|x;w) = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_3|x;w) = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

# Deep Neural Networks



$$z_i^{(k)} = g(\sum_j W_{i,j}^{(k-1,k)} z_j^{(k-1)})$$

**g = nonlinear activation function**

# Common Activation Functions

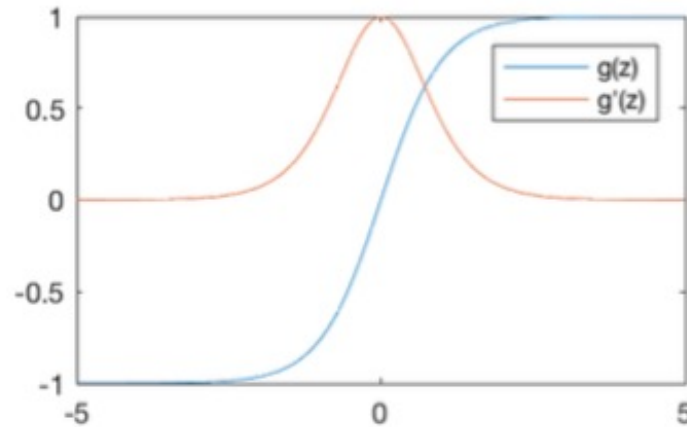## Sigmoid Function

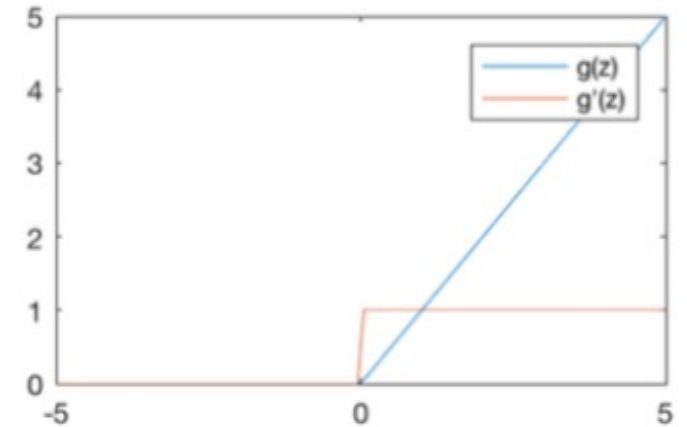$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$

## Hyperbolic Tangent

$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$

## Rectified Linear Unit (ReLU)

$$g(z) = \max(0, z)$$

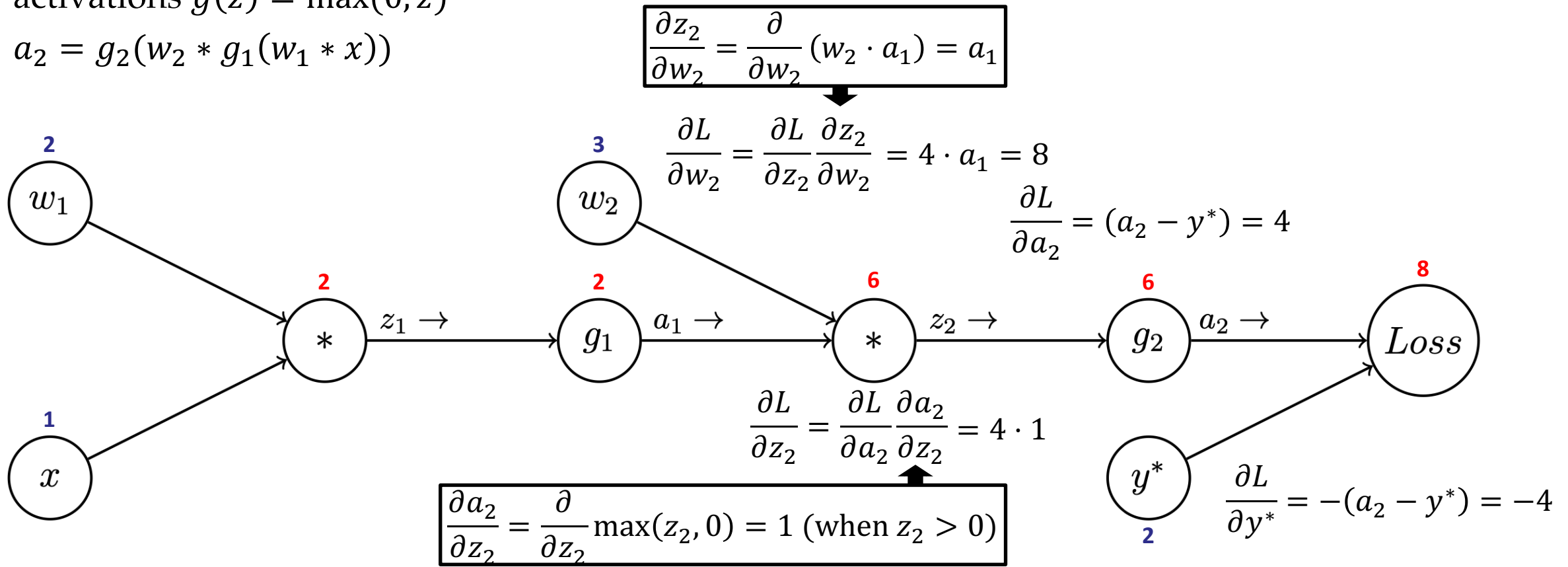$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

# Neural Network Properties

o Theorem (Universal Function Approximators). A two-layer neural network with a sufficient number of neurons can approximate any continuous function to any desired accuracy.

o Practical considerations:

    o Must have nonlinear activation function

    o Requires an arbitrarily large number of neurons:

        o Danger for overfitting (hence early stopping!)

        o No guarantee that we can do this on real-world compute

    o Often more efficient in practice to have more layers, less neurons

# Example: Automatic Differentiation

o Build a computation graph and apply chain rule: $f(x) = g(h(x))$ $\qquad$ $f'(x) = h'(x) \cdot g'(h(x))$

o Example: neural network with quadratic loss: $L(a_2, y^*) = \frac{1}{2}(a_2 - y^*)^2$ and ReLU activations $g(z) = \max(0, z)$
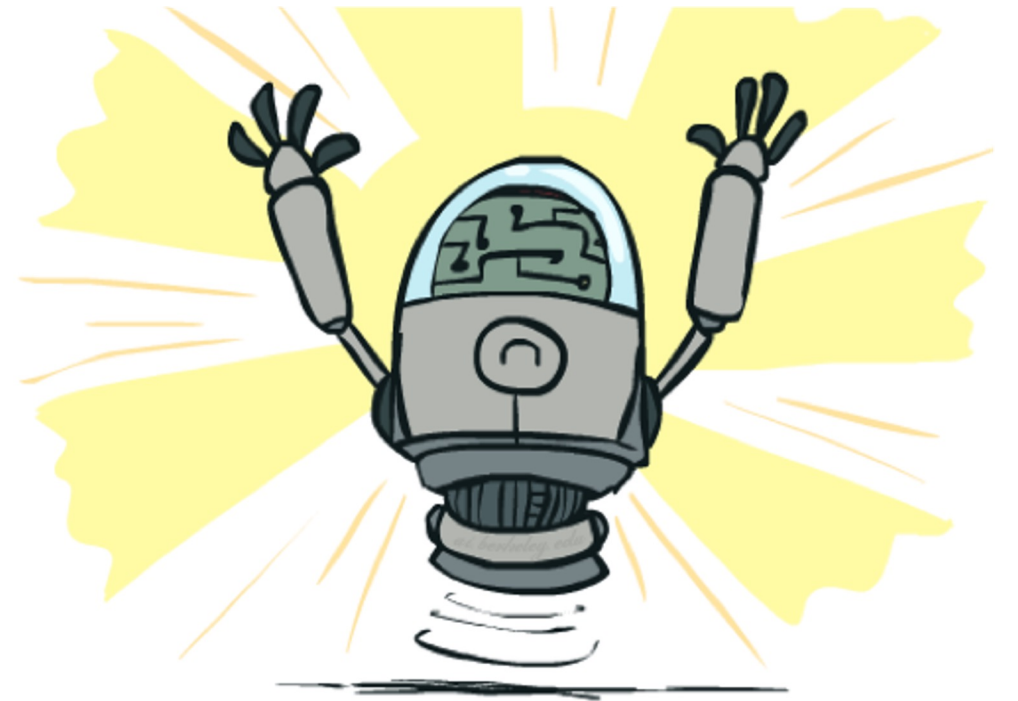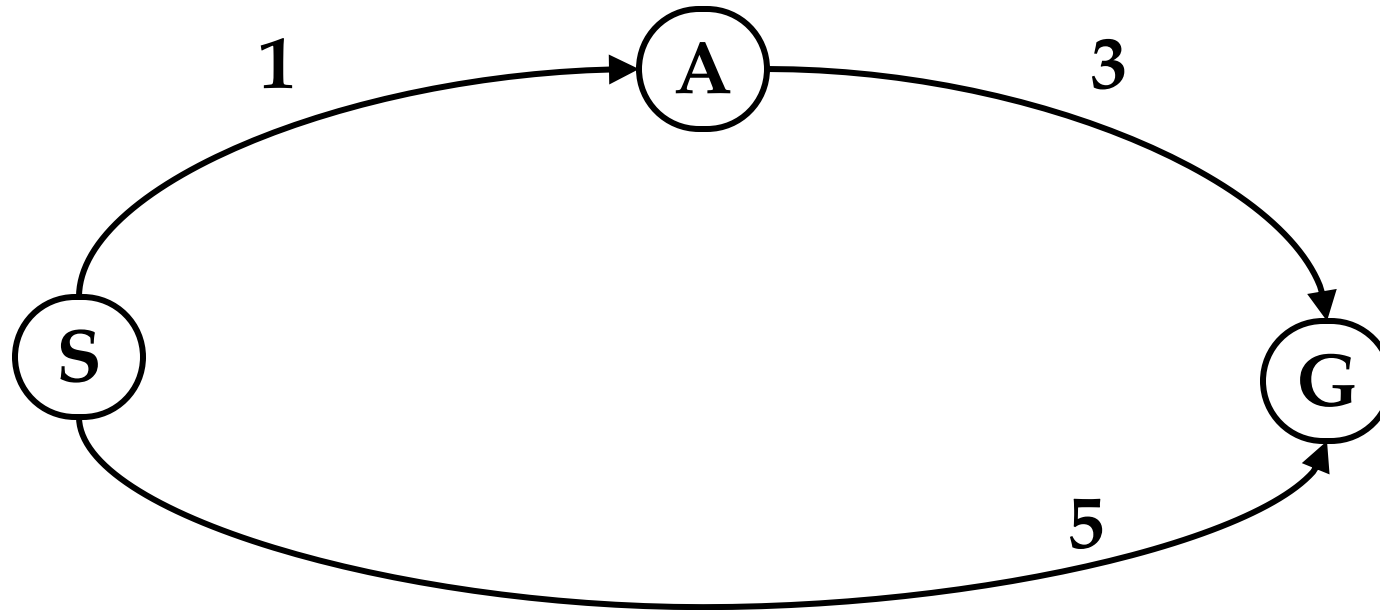
o $a_2 = g_2(w_2 * g_1(w_1 * x))$

$$\boxed{\frac{\partial z_2}{\partial w_2} = \frac{\partial}{\partial w_2}(w_2 \cdot a_1) = a_1}$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial z_2}\frac{\partial z_2}{\partial w_2} = 4 \cdot a_1 = 8$$

$$\frac{\partial L}{\partial a_2} = (a_2 - y^*) = 4$$

$$\frac{\partial L}{\partial z_2} = \frac{\partial L}{\partial a_2}\frac{\partial a_2}{\partial z_2} = 4 \cdot 1$$

$$\boxed{\frac{\partial a_2}{\partial z_2} = \frac{\partial}{\partial z_2}\max(z_2, 0) = 1 \text{ (when } z_2 > 0)}$$

$$\frac{\partial L}{\partial y^*} = -(a_2 - y^*) = -4$$

# Search

# A* Search

o Expand nodes based on sum: backward cost + forward cost
  - o f(n) = g(n) + h(n)
  - o g(n): cost to get to node
  - o h(n): heuristic of future costs

o We ideally want heuristic functions that satisfy:
  - o Admissibility: underestimate true cost to the goal
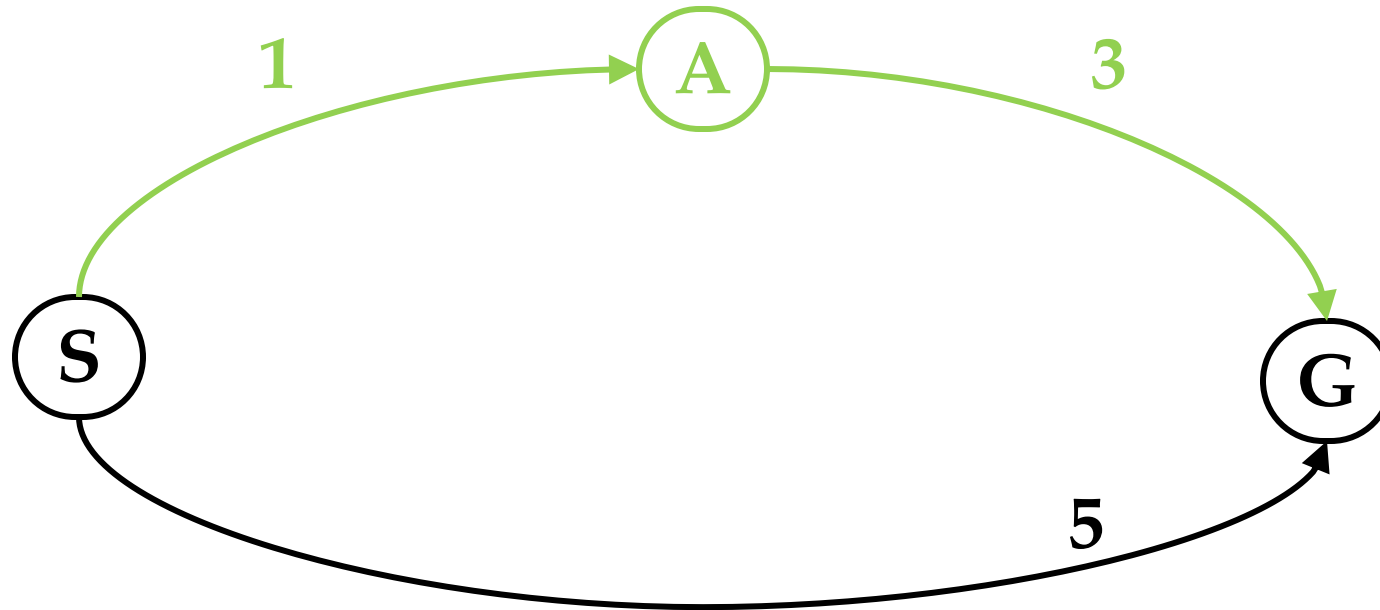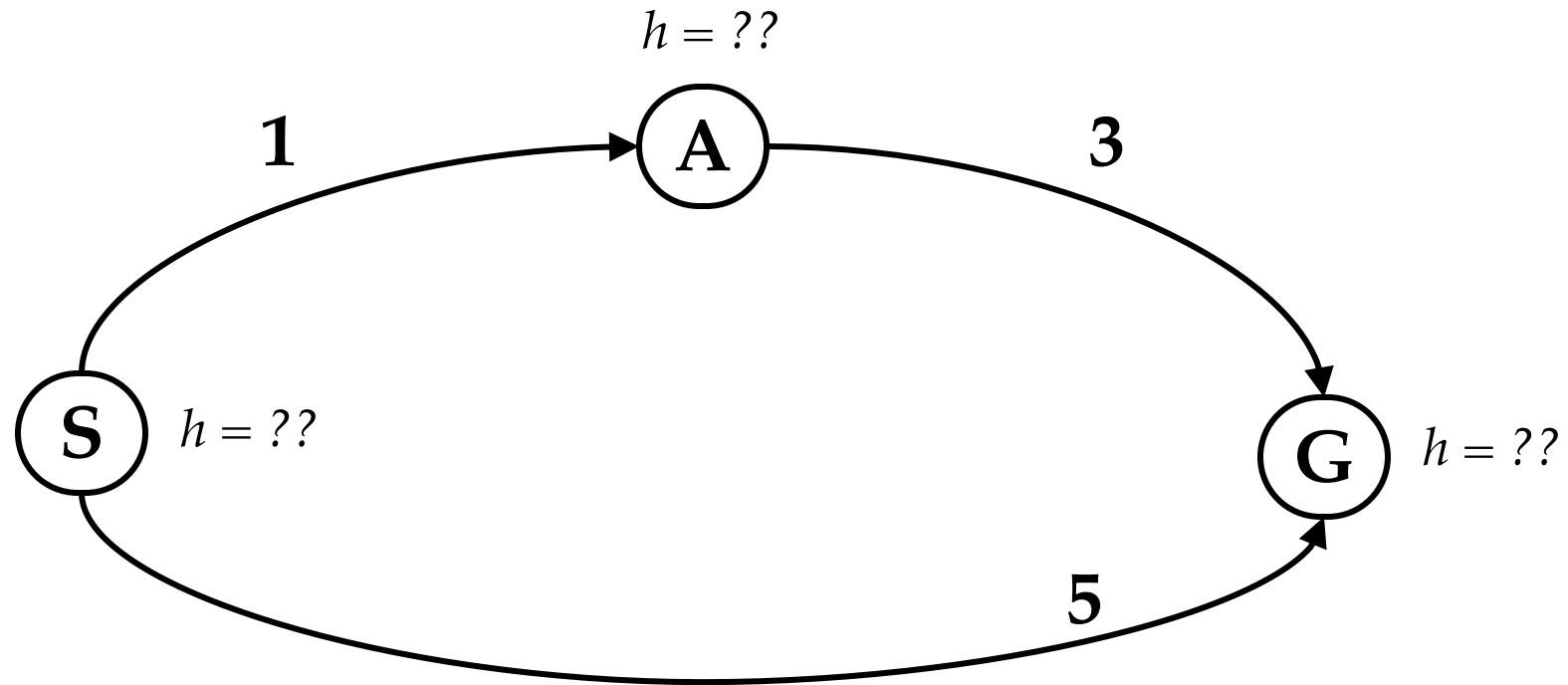  - o Consistency: "triangle inequality"

o Consistency => admissibility

# A* Search: Admissibility

# A* Search: Admissibility

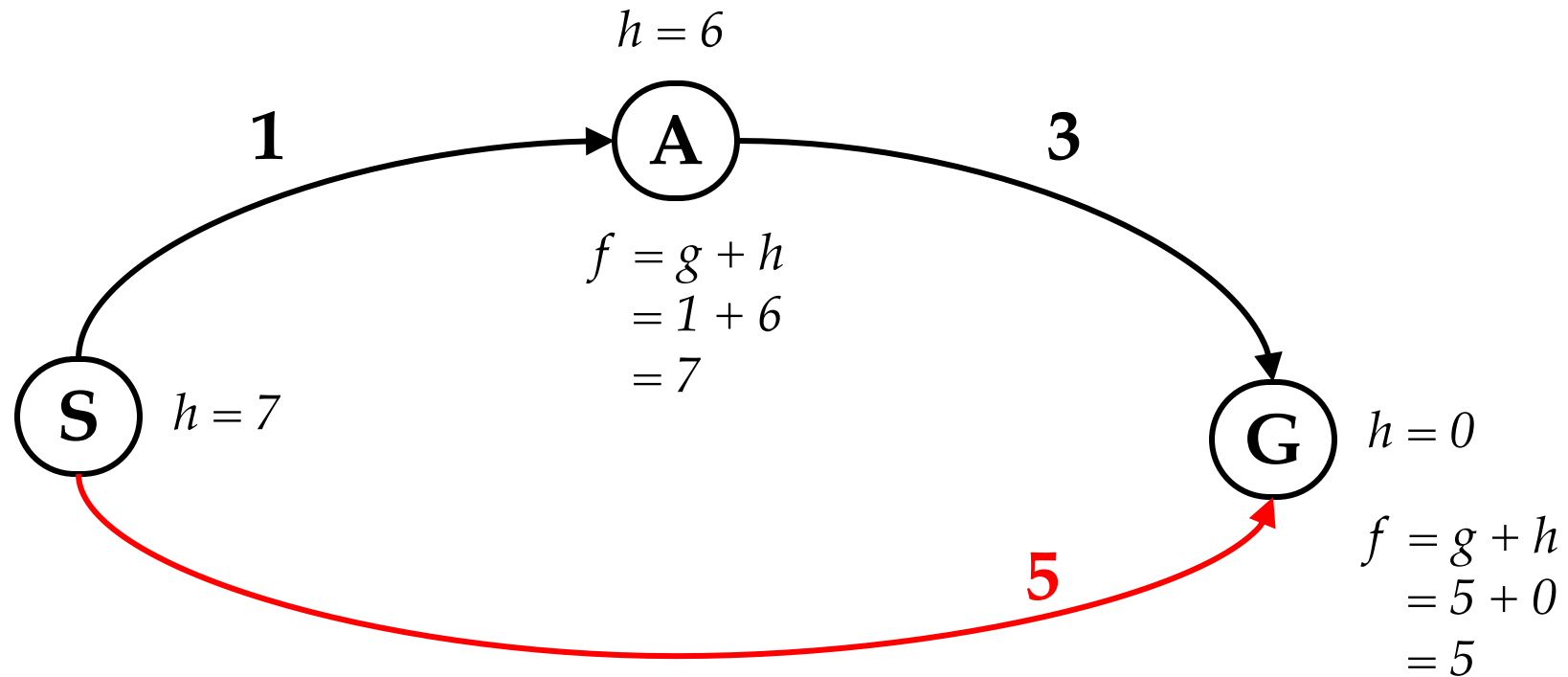Q: Where do heuristics come from?

A: We have to create them!



$h = 6$

1                    A                    3

$f = g + h$
$= 1 + 6$
$= 7$

S    $h = 7$

5

G    $h = 0$

$f = g + h$
$= 5 + 0$
$= 5$

Not the best heuristic…

# A* Search: Admissibility

Q: Where do heuristics come from?

A: We have to create them!

$$h = 6$$



$$f = g + h$$
$$= 1 + 6$$
$$= 7$$

$$h = 7$$

$$h = 0$$

$$f = g + h$$
$$= 5 + 0$$
$$= 5$$

Not the best heuristic…

Q: Where do heuristics come from?

A: We have to create them!



$h = 3$

**A**

**1**

**3**

$f = g + h$
$= 1 + 3$
$= 4$

**S** $h = 1$

**G** $h = 0$

**5**

$f = g + h$
$= 5 + 0$
$= 5$

What's a better heuristic?

# A* Search: Admissibility

Q: Where do heuristics come from?

A: We have to create them!



$h = 3$

**A**

$f = g + h$
$= 1 + 3$
$= 4$

**S** $h = 1$

**G** $h = 0$

1

3

5

$f = g + h$
$= 5 + 0$
$= 5$

What's a better heuristic?

**Admissible = Underestimates Cost to the Goal**
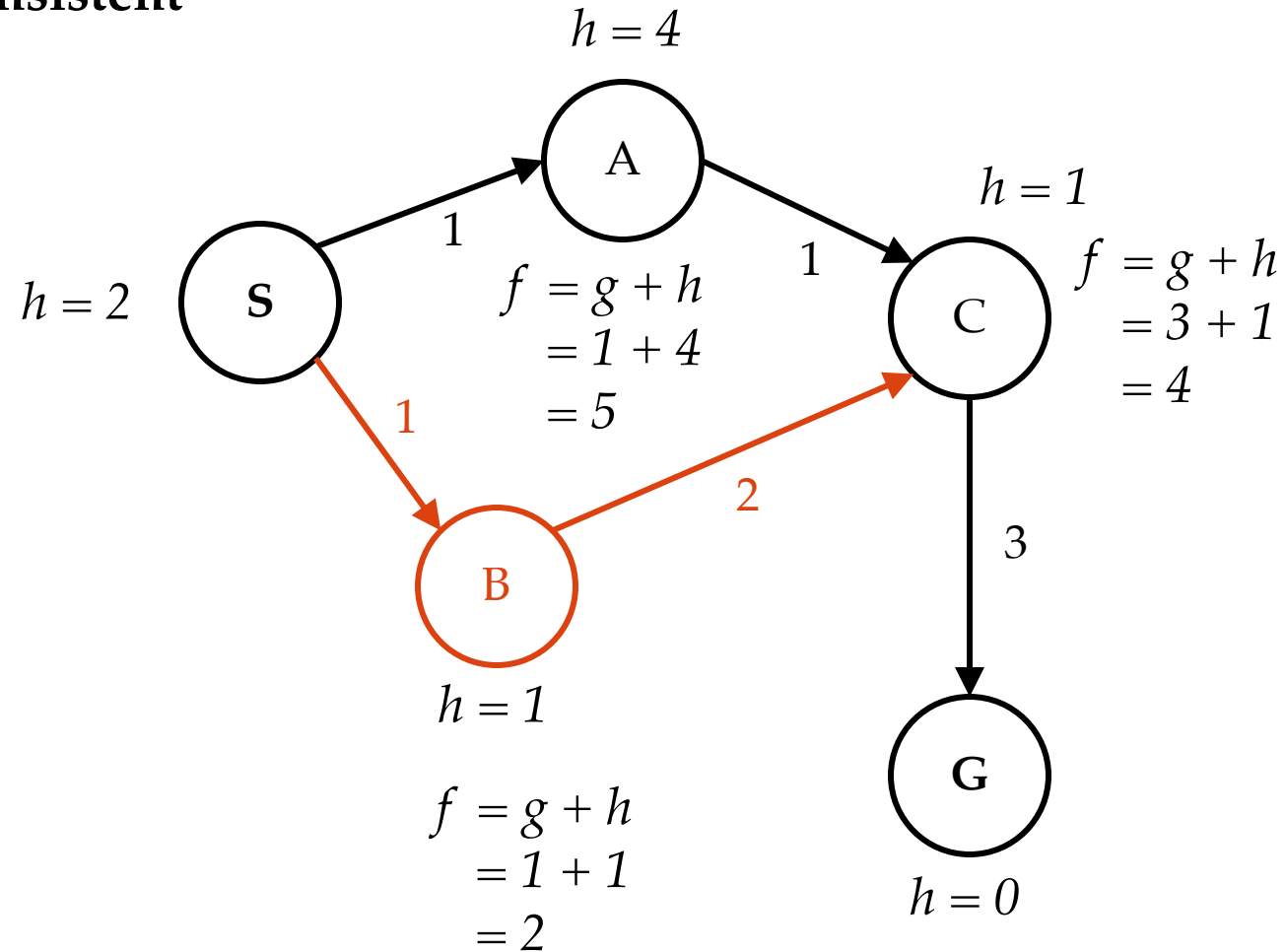
# A* Search: Consistency

# A* Search: Consistency

# A* Search: Consistency

$h = 4$

A

$h = 1$

$h = 2$   S

1

1

$f = g + h$
$= 1 + 4$
$= 5$

1

C

$f = g + h$
$= 3 + 1$
$= 4$

1

2

3

B

$h = 1$

$f = g + h$
$= 1 + 1$
$= 2$

G

$h = 0$

# A* Search: Consistency

This heuristic isn't **consistent**

$h = 4$

$h = 2$

**S**

A

1

$f = g + h$
$= 1 + 4$
$= 5$

$h = 1$

1

C

$f = g + h$
$= 3 + 1$
$= 4$

1

"Triangle inequality"
h(u) ≤ d(u,v) + h(v)

1

2

B

3

$h = 1$

$f = g + h$
$= 1 + 1$
$= 2$

**G**

$h = 0$

# A* Search: Consistency

This heuristic isn't **consistent**

$h = 4$

A

$h = 1$

$h = 2$  S

$f = g + h$
$= 1 + 4$
$= 5$

1

1

$f = g + h$
$= 3 + 1$
$= 4$

C

1

1

2

3

B

G

"Triangle inequality"
$h(u) \leq d(u,v) + h(v)$

Q: Is $h(A) \leq d(A,C) + h(C)$?

$h = 1$

$f = g + h$
$= 1 + 1$
$= 2$

$h = 0$

# A* Search: Consistency

This heuristic isn't **consistent**

$h = 4$

$h = 2$  **S**

A

$h = 1$

1

1

C

$f = g + h$
$= 1 + 4$
$= 5$

$f = g + h$
$= 3 + 1$
$= 4$

1

2

3

B

$h = 1$

**G**

"Triangle inequality"
$h(u) \leq d(u,v) + h(v)$

Q: Is $h(A) \leq d(A,C) + h(C)$?

A: No: $4 \nleq 1 + 1$

$f = g + h$
$= 1 + 1$
$= 2$

$h = 0$

# Summary of A*

o Tree search:
   o A* is optimal if heuristic is admissible
   o UCS is a special case ($h = 0$)

o Graph search:
   o A* optimal if heuristic is consistent
   o UCS optimal ($h = 0$ is consistent)

o Consistency implies admissibility

o In general, most natural admissible heuristics tend to be consistent, especially if it comes from a relaxed problem

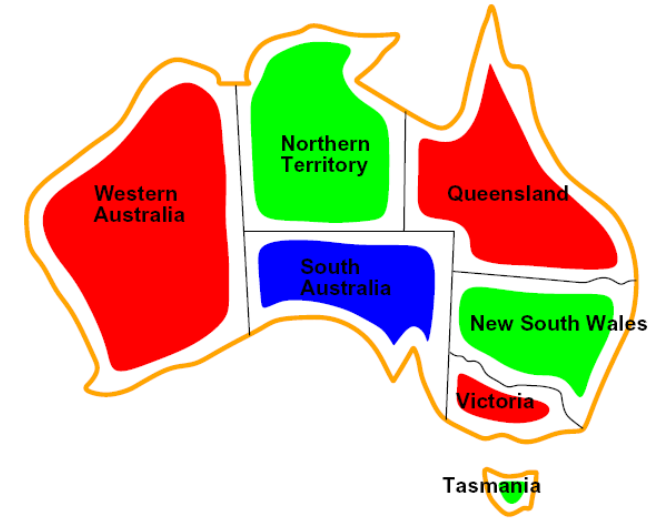# Constraint Satisfaction Problems

# Example: Map Coloring

o Variables: WA, NT, Q, NSW, V, SA, T

o Domains: $D = \{red, green, blue\}$

o Constraints: adjacent regions must have different colors

    Implicit: $WA \neq NT$

    Explicit: $(WA, NT) \in \{(red, green), (red, blue), \ldots\}$

o Solutions are assignments satisfying all constraints, e.g.:

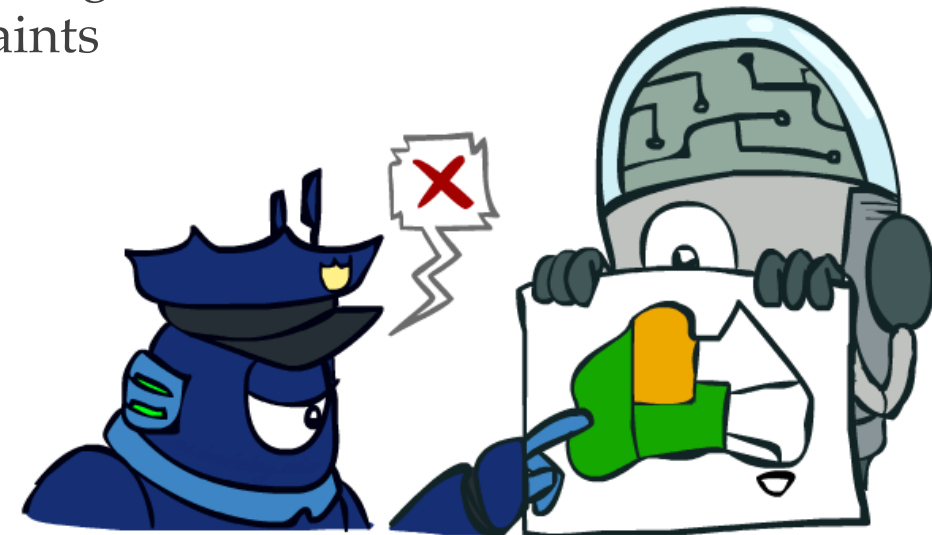    {WA=red, NT=green, Q=red, NSW=green, V=red, SA=blue, T=green}

# General Approach #1: Backtracking Search

- Backtracking search is the basic uninformed algorithm for solving CSPs

- Idea 1: One variable at a time
  - Variable assignments are commutative, so fix ordering -> better branching factor!
  - I.e., [WA = red then NT = green] same as [NT = green then WA = red]
  - Only need to consider assignments to a single variable at each step

- Idea 2: Check constraints as you go
  - I.e. consider only values which do not conflict previous assignments
  - Might have to do some computation to check the constraints
  - "Incremental goal test"

- Depth-first search with these two improvements is called *backtracking search* (not the best name)

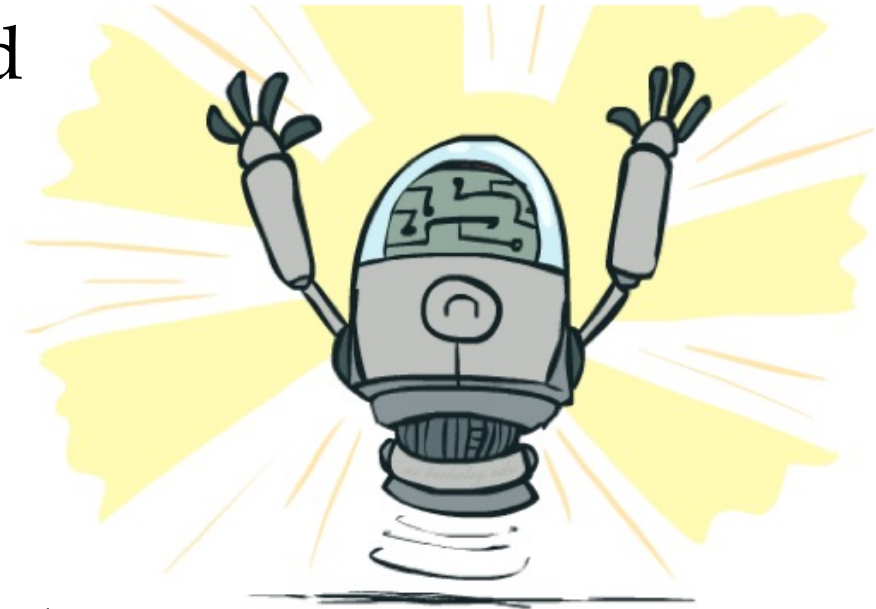- Can solve n-queens for n ≈ 25

# Improving Backtracking

General-purpose ideas give huge gains in speed

1. Ordering:
   o Which variable should be assigned next?
   o In what order should its values be tried?

2. Filtering: Can we detect inevitable failure early?

3. Leveraging the structure of the constraint graph

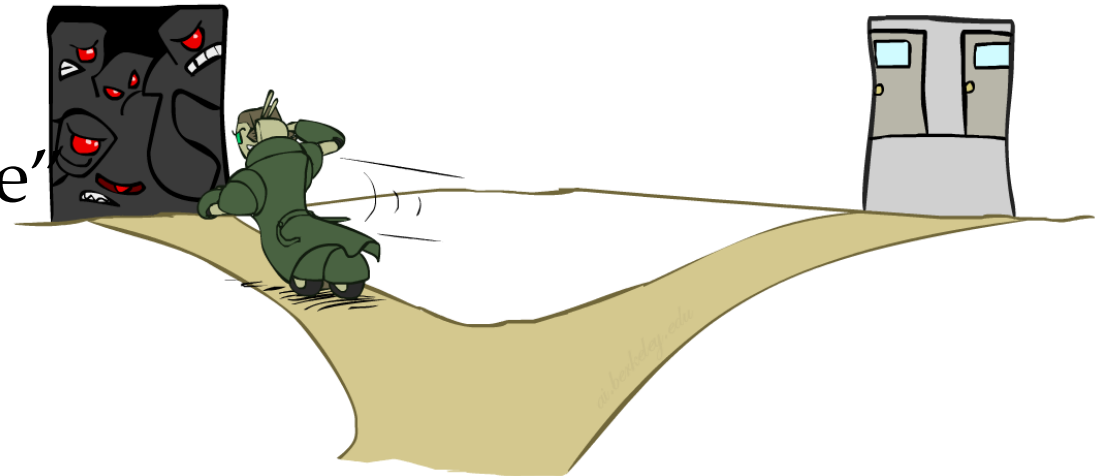# Ordering: Minimum Remaining Values

○ Variable Ordering: Minimum remaining values (MRV):

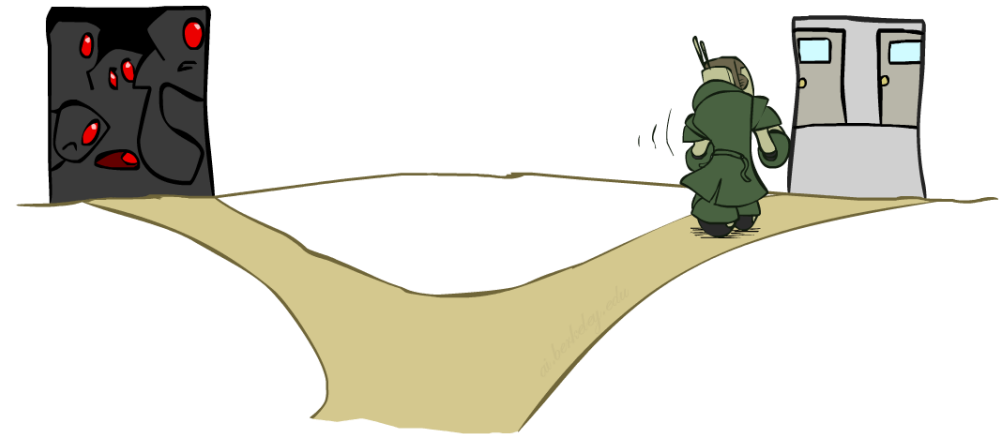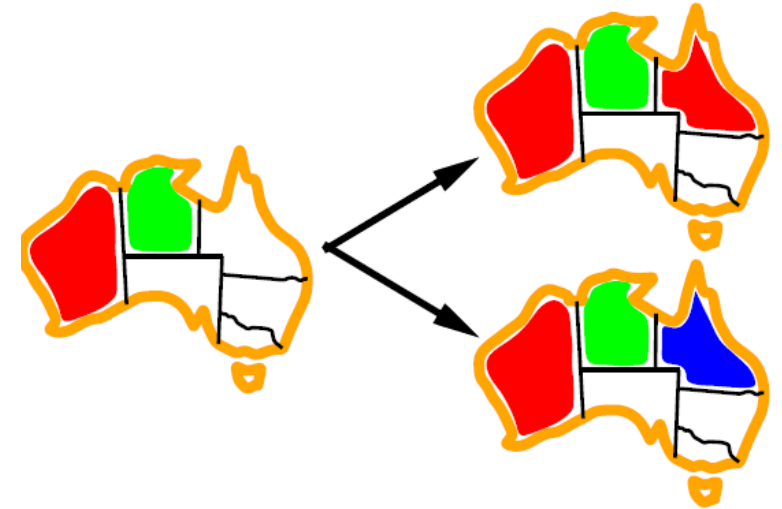   ○ Choose the variable with the fewest legal values left in its domain

○ Why min rather than max?

○ Also called "most constrained variable"

○ "Fail-fast" ordering
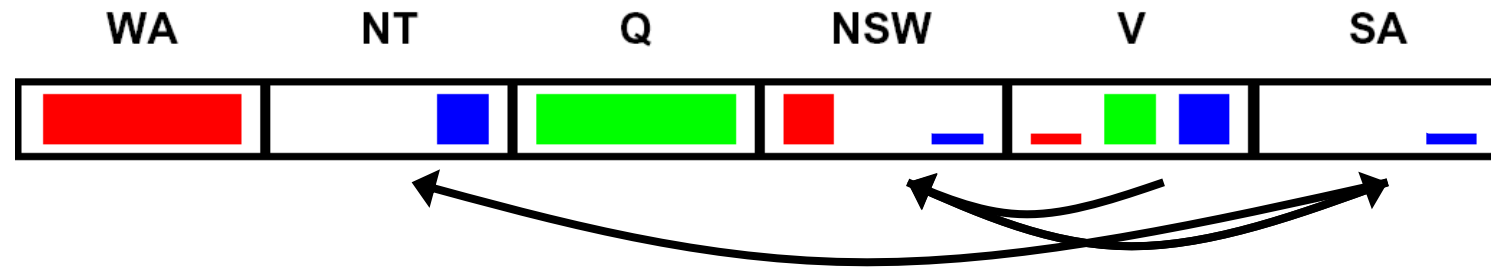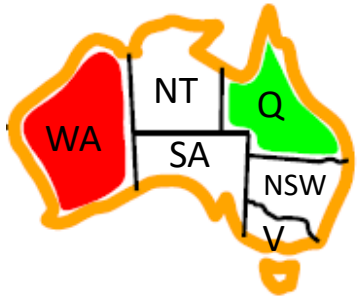
# Ordering: Least Constraining Value

o Value Ordering: Least Constraining Value

    o Given a choice of variable, choose the *least constraining value*

    o I.e., the one that rules out the fewest values in the remaining variables

    o Note that it may take some computation to determine this! (E.g., rerunning filtering)

o Why least rather than most?

o Combining these ordering ideas makes 1000 queens feasible

# Filtering: Arc Consistency

o A simple form of propagation makes sure *all* arcs are consistent:



o Important: If X loses a value, neighbors of X need to be rechecked!
o Arc consistency detects failure earlier than forward checking
o Can be run as a preprocessor or after each assignment
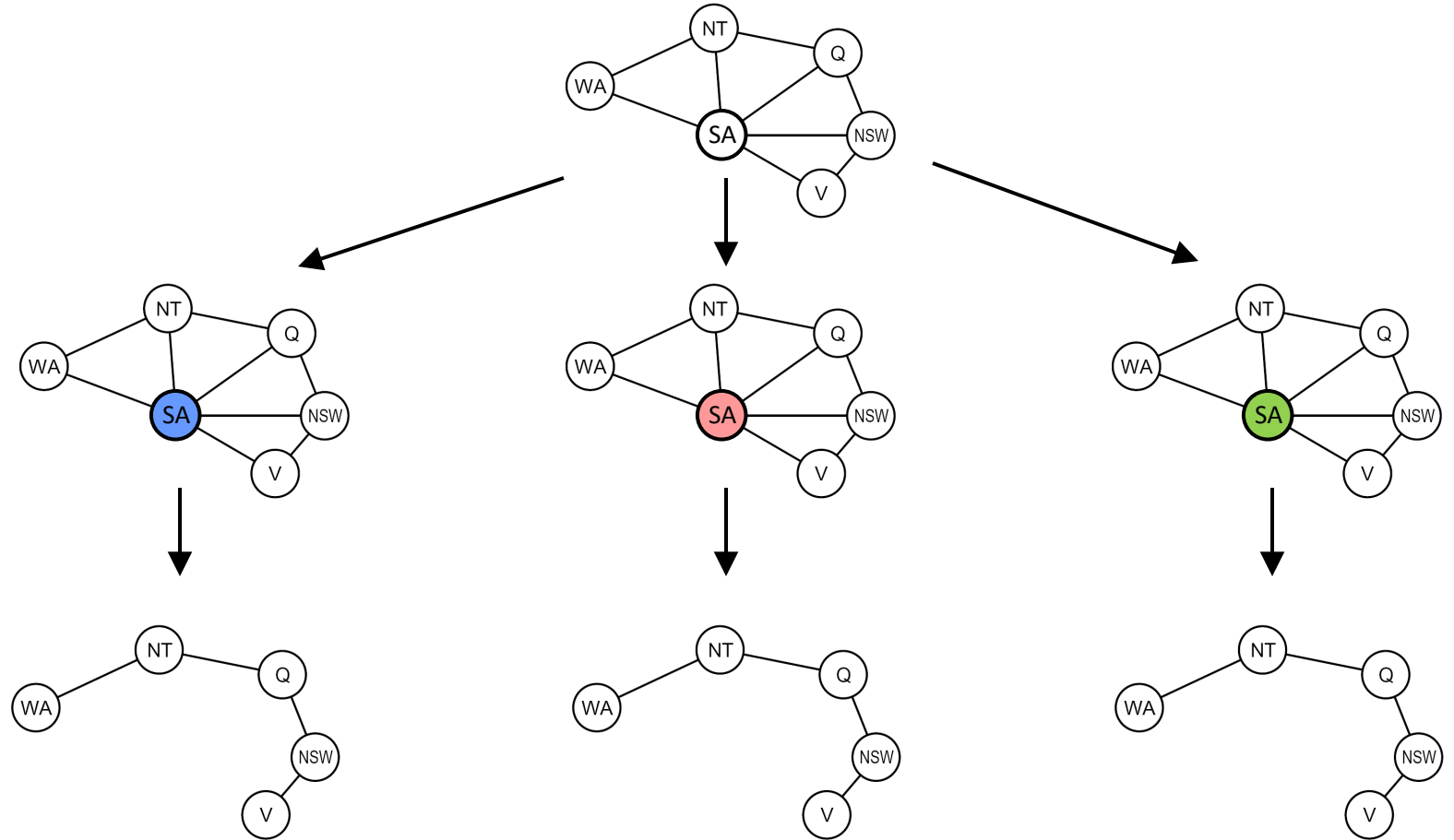
*Remember: Delete from the tail!*

# Leveraging Structure: Cutsets



Choose a cutset

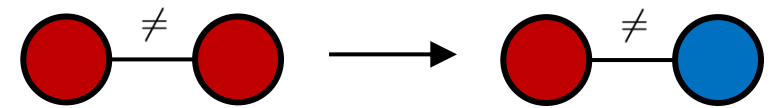Instantiate the cutset
(all possible ways)

Compute residual CSP
for each assignment

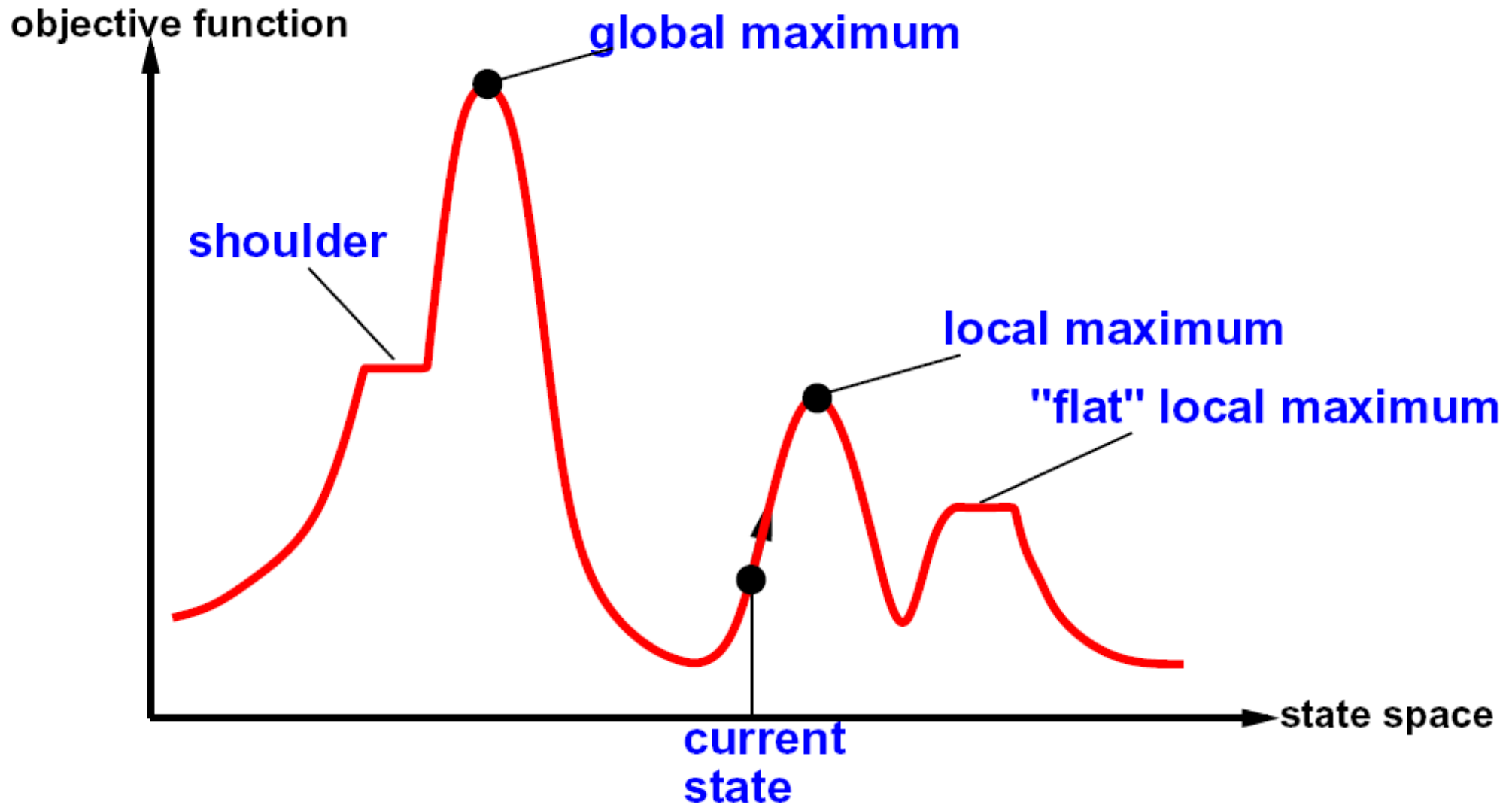Solve the residual CSPs
(tree structured)

# General Approach #2: Iterative Improvement

o Local search methods typically work with "complete" states, i.e., all variables assigned

o To apply to CSPs:
  o Take an assignment with unsatisfied constraints
  o Operators *reassign* variable values
  o No fringe!  Live on the edge.

o Algorithm: While not solved,
  o Variable selection: randomly select any conflicted variable
  o Value selection: min-conflicts heuristic:
    o Choose a value that violates the fewest constraints
    o I.e., hill climb with $h(x)$ = total number of violated constraints

# Hill Climbing Diagram

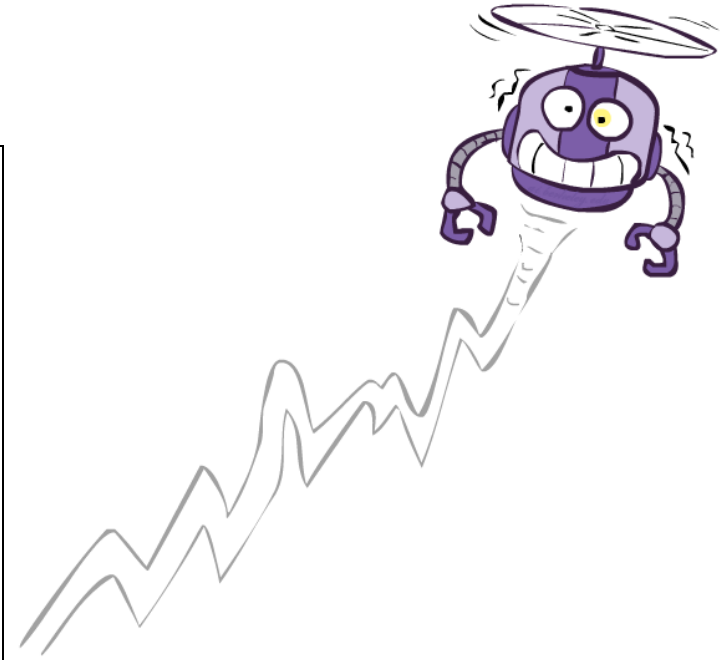# Simulated Annealing

o Idea:  Escape local maxima by allowing downhill moves
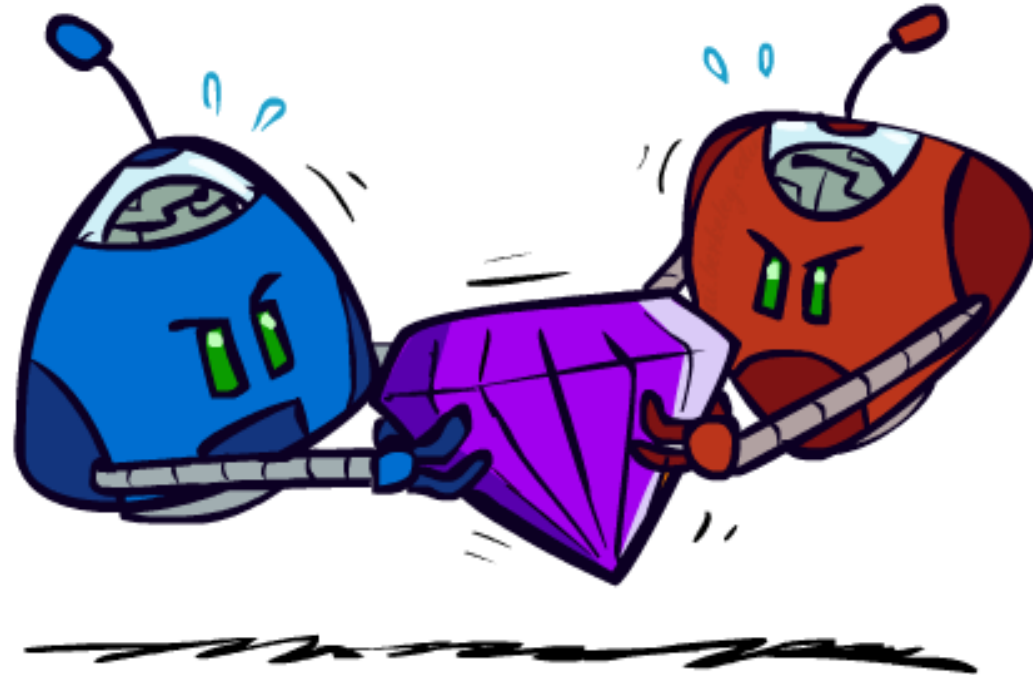  o But make them rarer as time goes on

**function** SIMULATED-ANNEALING( *problem, schedule* ) **returns** a solution state
   **inputs**: *problem*, a problem
           *schedule*, a mapping from time to "temperature"
   **local variables**: *current*, a node
              *next*, a node
              $T$, a "temperature" controlling prob. of downward steps

  *current* ← MAKE-NODE(INITIAL-STATE[*problem*])
  **for** $t$ ← 1 **to** ∞ **do**
    $T$ ← *schedule*[*t*]
    **if** $T = 0$ **then return** *current*
    *next* ← a randomly selected successor of *current*
    $\Delta E$ ← VALUE[*next*] − VALUE[*current*]
    **if** $\Delta E > 0$ **then** *current* ← *next*
    **else** *current* ← *next* only with probability $e^{\Delta E/T}$

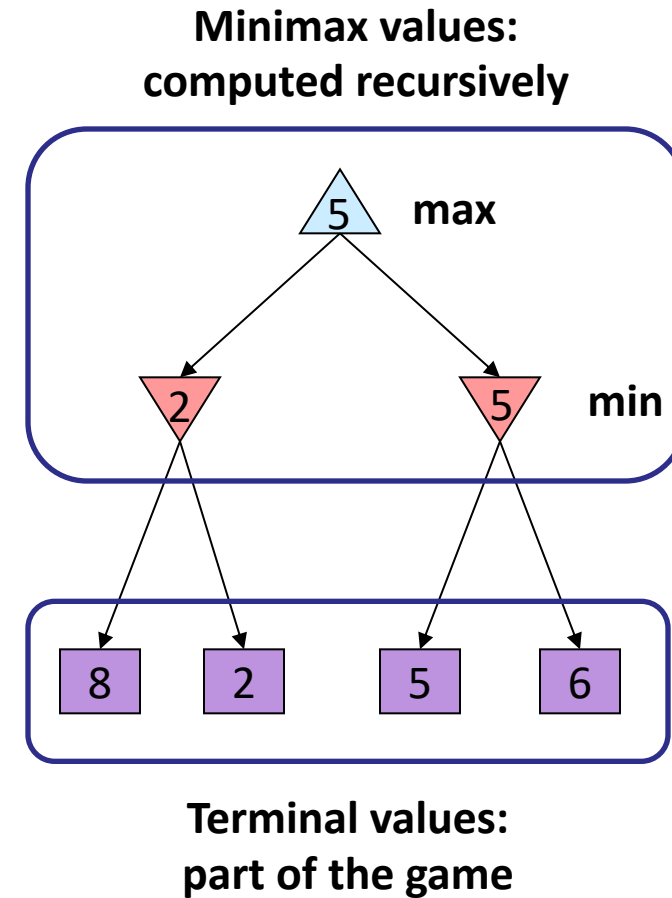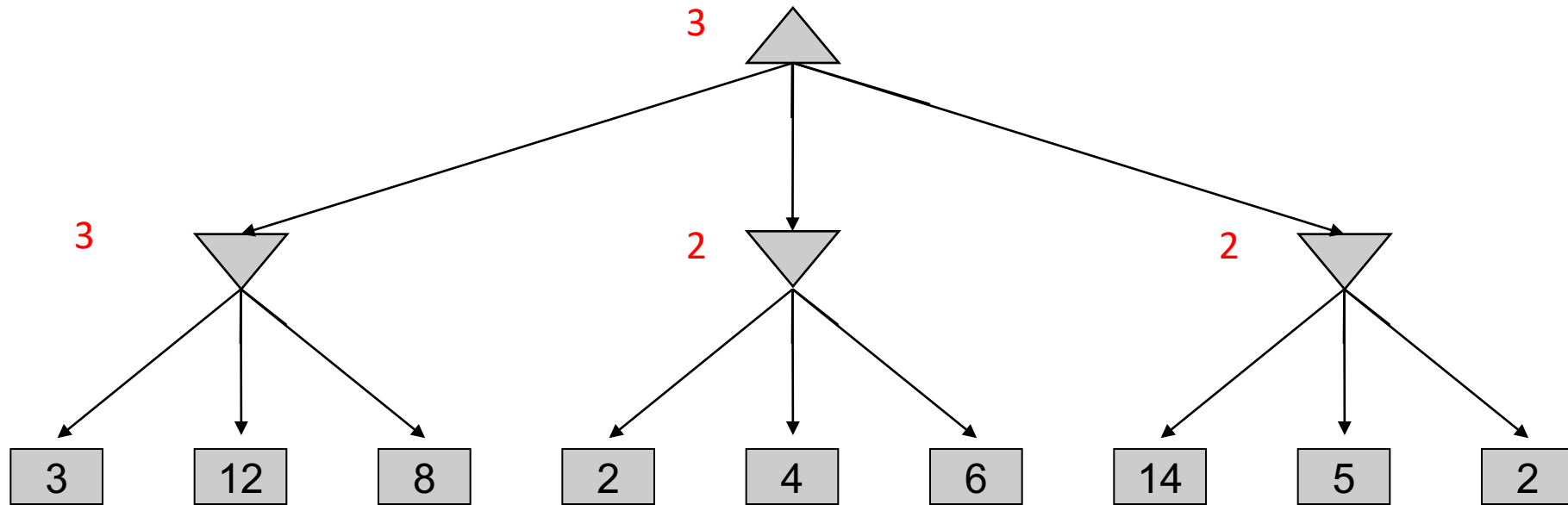# Game Trees

# Adversarial Search (Minimax)

o Deterministic, zero-sum games:
  o Tic-tac-toe, chess, checkers
  o One player maximizes result
  o The other minimizes result

o Minimax search:
  o A state-space search tree
  o Players alternate turns
  o Compute each node's <span style="color:red">minimax value:</span> the best achievable utility against a rational (optimal) adversary
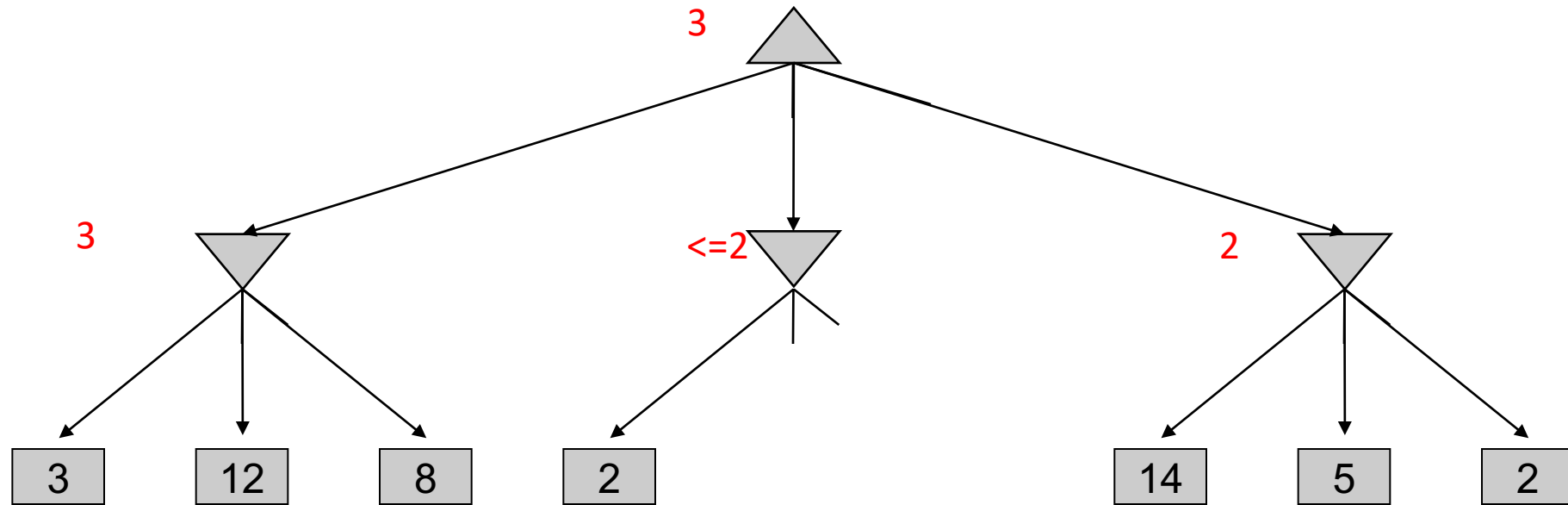


**Minimax values:**
**computed recursively**

max

min

**Terminal values:**
**part of the game**

# Minimax Example
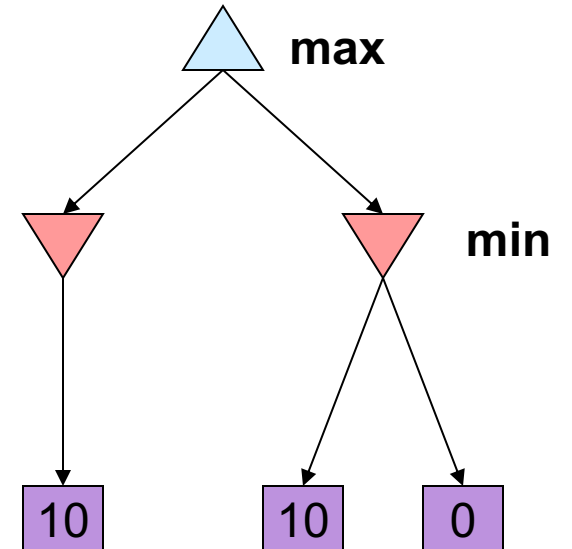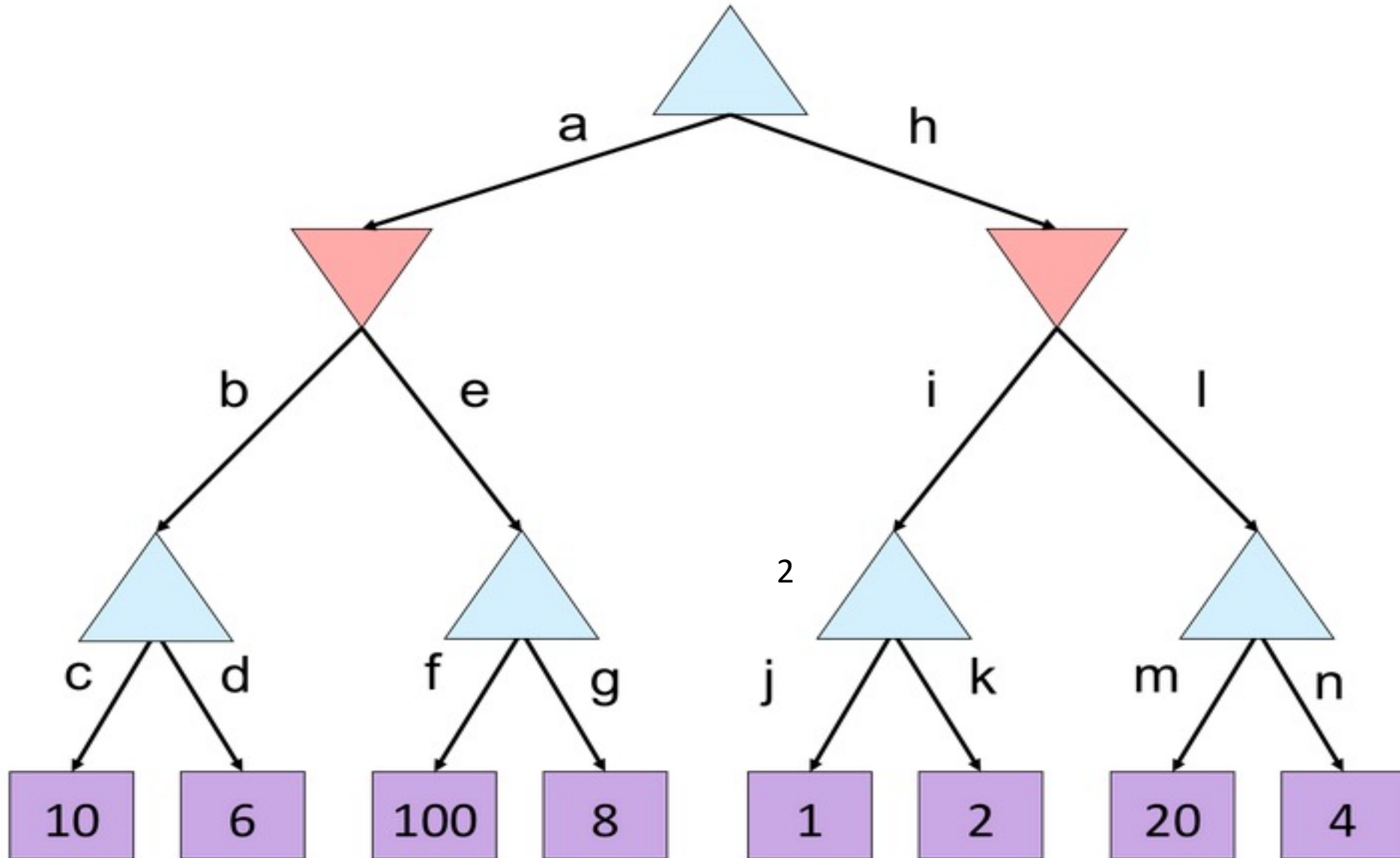
# Minimax Example: Pruning

# Alpha-Beta Pruning Properties

o This pruning has <span style="color:red">no effect</span> on minimax value computed for the root!

o Values of intermediate nodes might be wrong
  o Important: children of the root may have the wrong value
  o So the most naïve version won't let you do action selection

o Good child ordering improves effectiveness of pruning

o With "perfect ordering":
  o Time complexity drops to $O(b^{m/2})$
  o Doubles solvable depth!
  o Full search of, e.g. chess, is still hopeless…

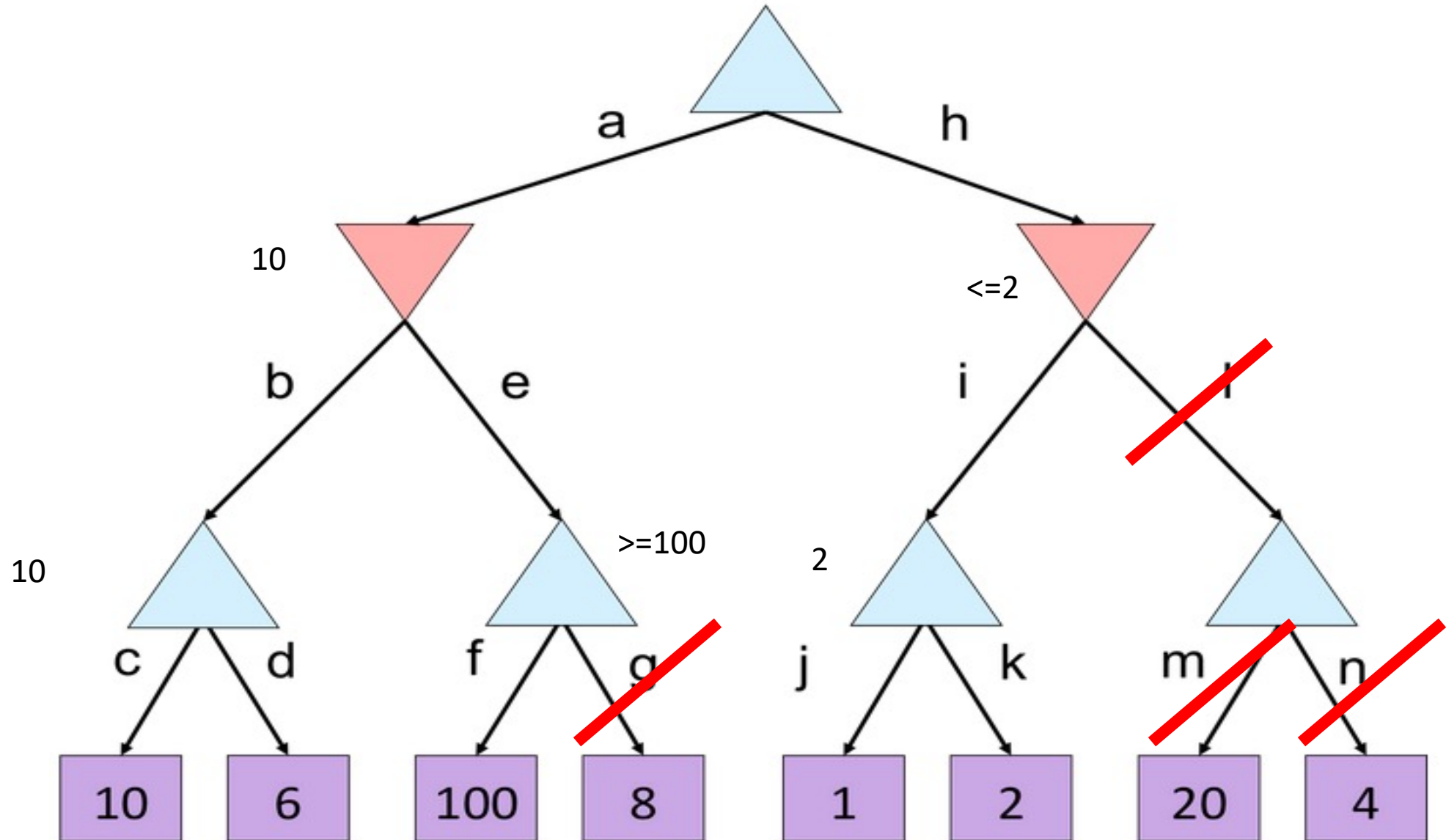o This is a simple example of <span style="color:red">metareasoning</span> (computing about what to compute)
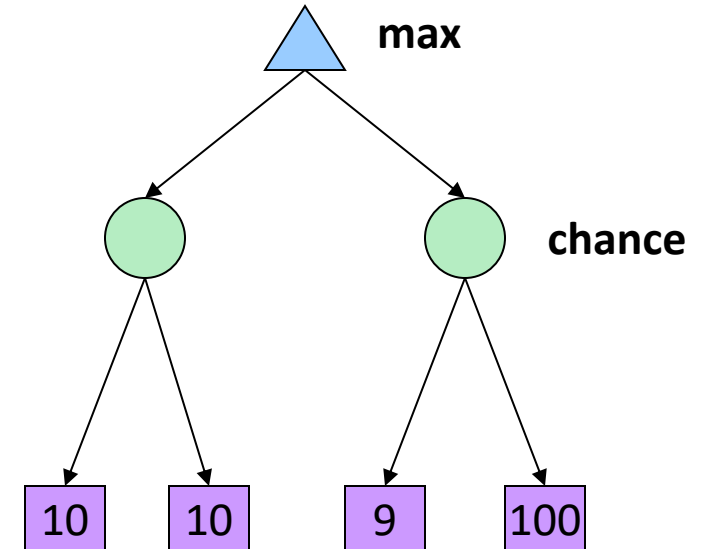
**max**

**min**

10    10    0

# Alpha-Beta Quiz 2

# Expectimax Search

o Why wouldn't we know what the result of an action will be?
   o Explicit randomness: rolling dice
   o Unpredictable opponents: the ghosts respond randomly
   o Unpredictable humans: humans are not perfect
   o Actions can fail: when moving a robot, wheels might slip

o Values should now reflect average-case (expectimax) outcomes, not worst-case (minimax) outcomes

o Expectimax search: compute the average score under optimal play
   o Max nodes as in minimax search
   o Chance nodes are like min nodes but the outcome is uncertain
   o Calculate their expected utilities
   o I.e. take weighted average (expectation) of children

# Remaining Topics

Bayes Nets:
- Inference by enumeration
- Variable elimination
- D-separation
- Sampling approaches

HMMs:
- Forward algorithm
- Viterbi algorithm
- Particle filtering

Decision networks and VPIs

Out of scope: learning theory, decision tree classifiers, details of non-SGD optimizers (e.g., NAG, Adagrad, Adam), NLP/CV/RL