

Per-Decision Option Discounting

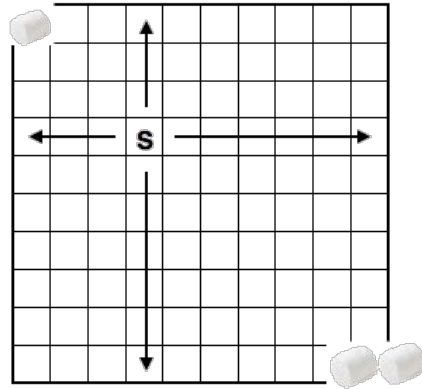
Anna Harutyunyan, Peter Vrancx,
Philippe Hamel, Ann Nowe,
Doina Precup

Motivation: Agents that reason over long temporal horizons

Motivation: Agents that reason over long temporal horizons

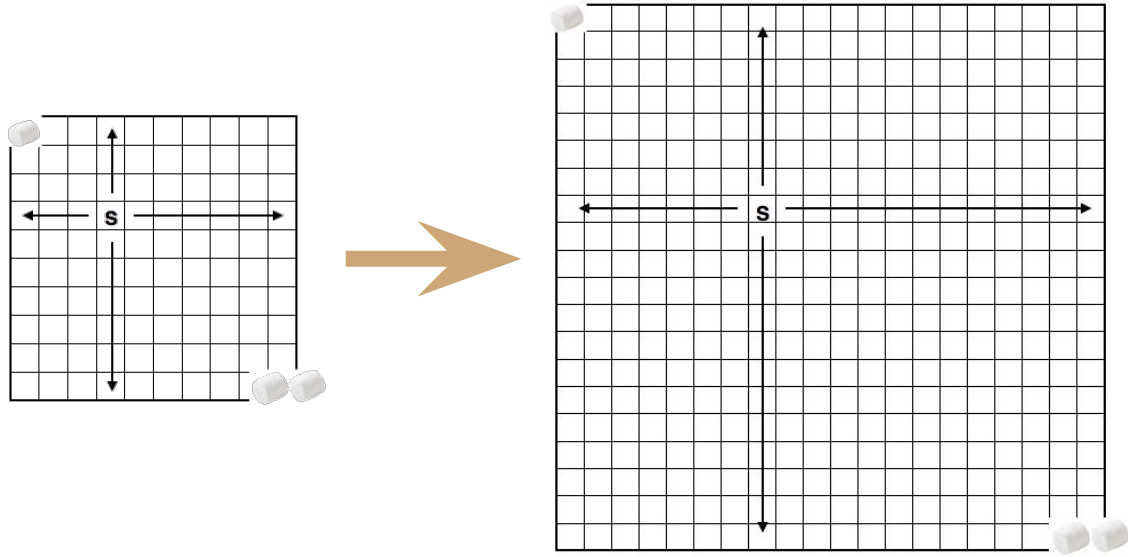
Horizon depends on discount γ

Motivation: Agents that reason over long temporal horizons



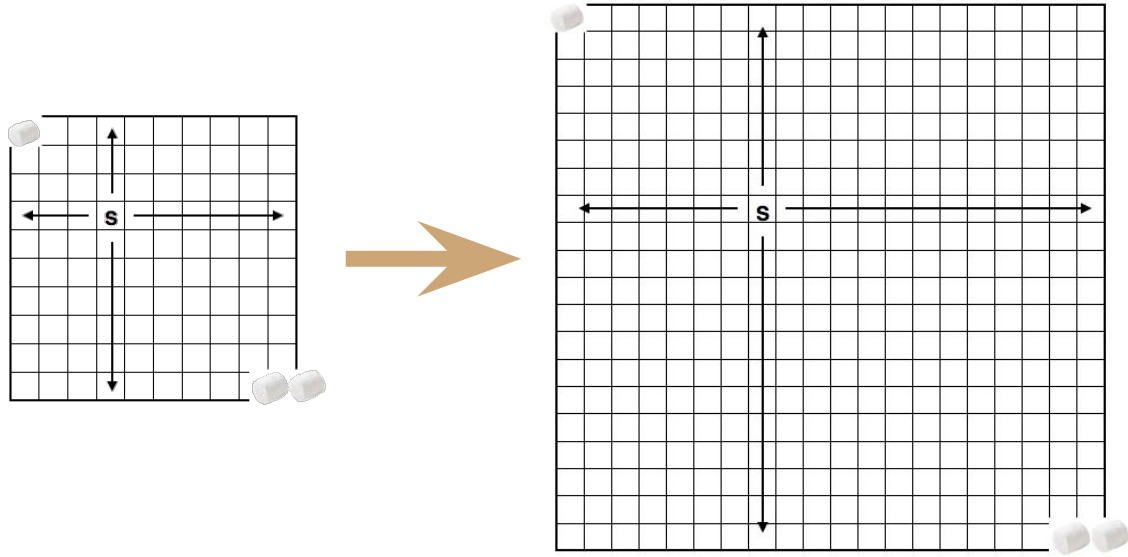
Horizon depends on discount γ

Motivation: Agents that reason over long temporal horizons



Horizon depends on discount γ
Larger grid requires a larger γ

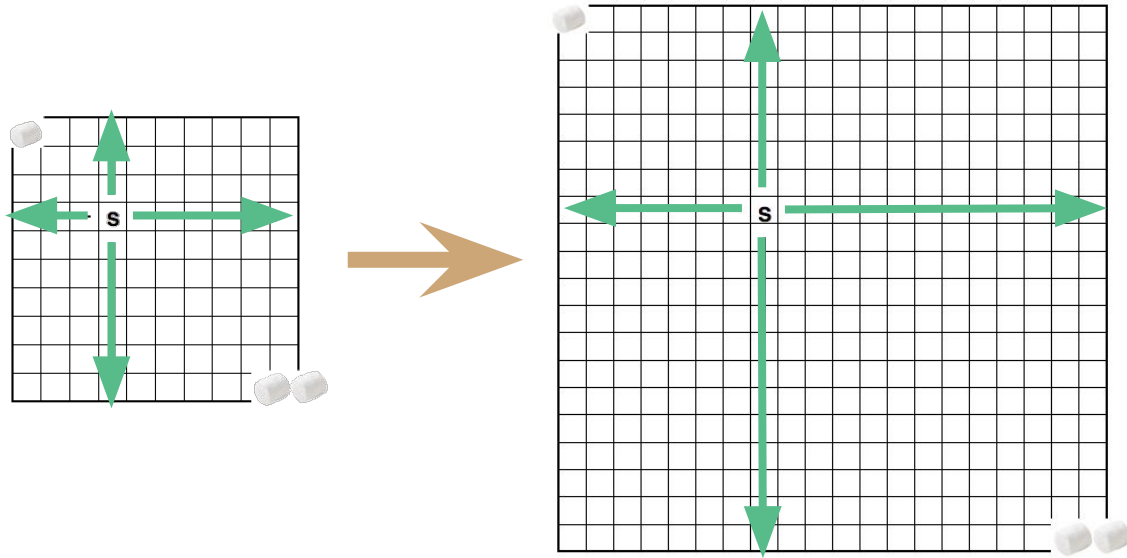
Motivation: Agents that reason over long temporal horizons



Horizon depends on discount γ
Larger grid requires a larger γ

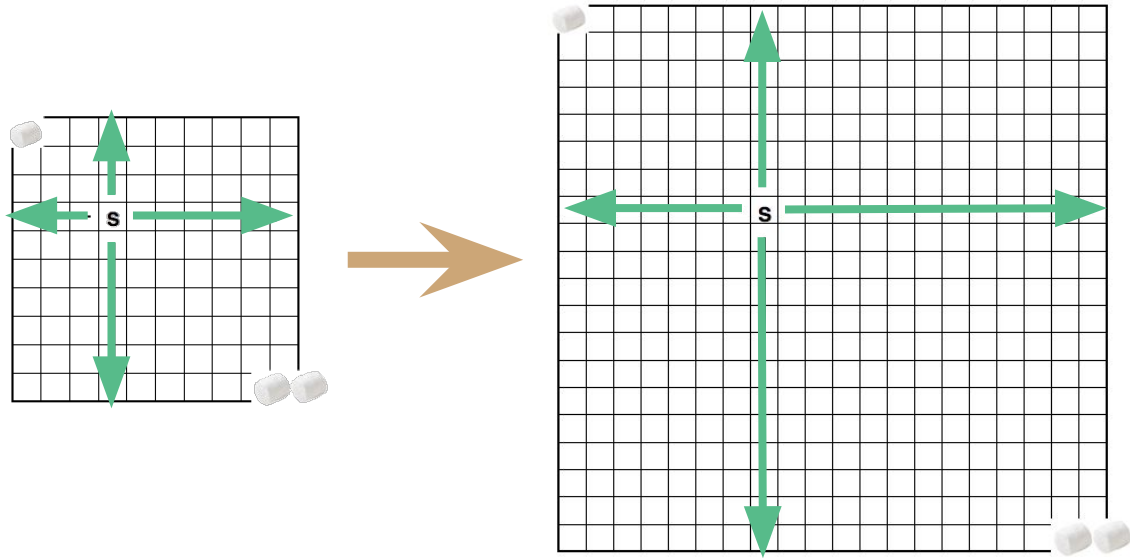
Large γ -s are inefficient in practice :(

Motivation: Agents that reason over long temporal horizons



Horizon depends on discount γ
Larger grid requires a larger γ
Temporal abstraction?

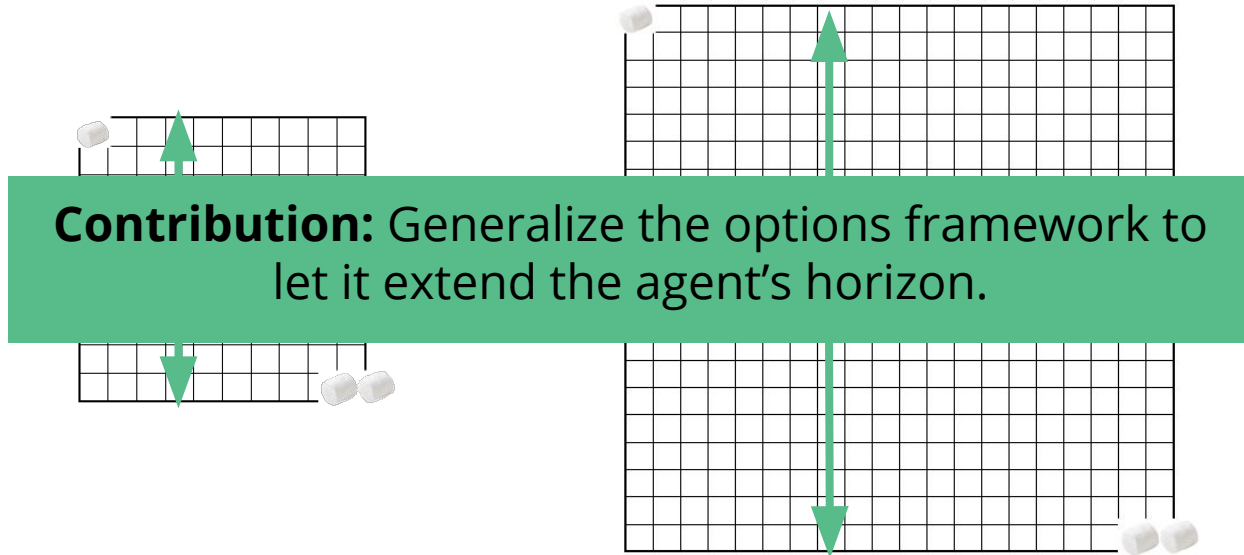
Motivation: Agents that reason over long temporal horizons



Horizon depends on discount γ
Larger grid requires a larger γ

Temporal abstraction? Options still tied to γ !

Motivation: Agents that reason over long temporal horizons



Horizon depends on discount γ
Larger grid requires a larger γ

Temporal abstraction? Options still tied to γ !

The Options Framework

Reward model:

$$R^o(s) \stackrel{\text{def}}{=} \mathbb{E}_{D:s|o} \left[\sum_{t=0}^{D-1} \gamma^t R_{t+1} \right]$$

Transition model:

$$P^o(s'|s) \stackrel{\text{def}}{=} \mathbb{E}_{D:s \rightarrow s'|o} \left[\gamma^D \mathbb{I}_{S_D=s'} \right]$$

The Options Framework

Reward model:

$$R^o(s) \stackrel{\text{def}}{=} \mathbb{E}_{D:s|o} \left[\sum_{t=0}^{D-1} \gamma^t R_{t+1} \right]$$

Transition model:

$$P^o(s'|s) \stackrel{\text{def}}{=} \mathbb{E}_{D:s \rightarrow s'|o} \left[\gamma^D \mathbb{1}_{S_D=s'} \right]$$

Options with Time Dilation

Reward model:

$$R^o(s) \stackrel{\text{def}}{=} \mathbb{E}_{D:s|o} \left[\sum_{t=0}^{D-1} \gamma^t R_{t+1} \right]$$



$$R^o(s) \stackrel{\text{def}}{=} \mathbb{E}_{D:s|o} \left[\sum_{t=0}^{D-1} \gamma_r^t R_{t+1} \right]$$

Transition model:

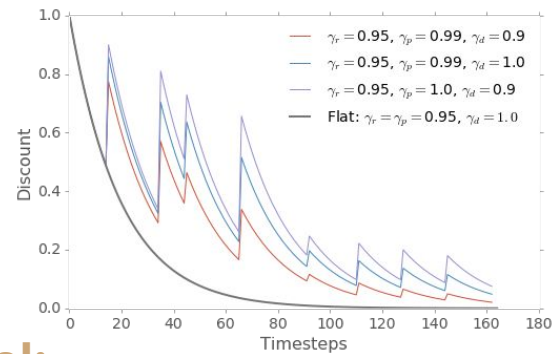
$$P^o(s'|s) \stackrel{\text{def}}{=} \mathbb{E}_{D:s \rightarrow s'|o} \left[\gamma^D \mathbb{1}_{S_D=s'} \right]$$



(2) per-decision (1) decouple

$$P^o(s'|s) \stackrel{\text{def}}{=} \gamma_d \mathbb{E}_{D:s \rightarrow s'|o} \left[\gamma_p^D \mathbb{1}_{S_D=s'} \right]$$

Options with Time Dilation



Reward model:

$$R^o(s) \stackrel{\text{def}}{=} \mathbb{E}_{D:s|o} \left[\sum_{t=0}^{D-1} \gamma^t R_{t+1} \right]$$

↓

$$R^o(s) \stackrel{\text{def}}{=} \mathbb{E}_{D:s|o} \left[\sum_{t=0}^{D-1} \gamma_r^t R_{t+1} \right]$$

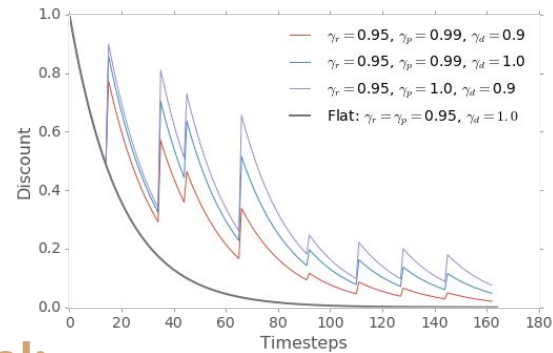
Transition model:

$$P^o(s'|s) \stackrel{\text{def}}{=} \mathbb{E}_{D:s \rightarrow s'|o} \left[\gamma^D \mathbb{1}_{S_D=s'} \right]$$

↓ (2) per-decision (1) decouple

$$P^o(s'|s) \stackrel{\text{def}}{=} \mathbb{E}_{D:s \rightarrow s'|o} \left[\gamma_d \mathbb{1}_{S_D=s'} \right]$$

Options with Time Dilation



Reward model:

$$R^o(s) \stackrel{\text{def}}{=} \mathbb{E}_{D:s|o} \left[\sum_{t=0}^{D-1} \gamma^t R_{t+1} \right]$$

↓

$$R^o(s) \stackrel{\text{def}}{=} \mathbb{E}_{D:s|o} \left[\sum_{t=0}^{D-1} \gamma_r^t R_{t+1} \right]$$

Transition model:

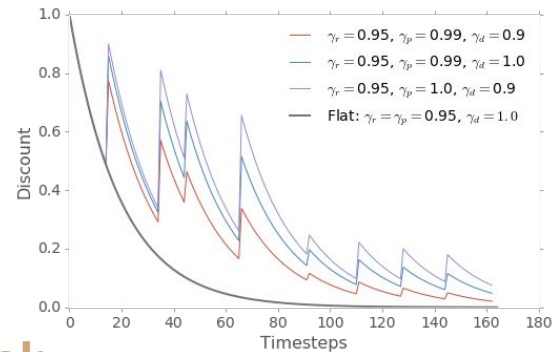
$$P^o(s'|s) \stackrel{\text{def}}{=} \mathbb{E}_{D:s \rightarrow s'|o} \left[\gamma^D \mathbb{1}_{S_D=s'} \right]$$

↓ (2) per-decision (1) decouple

$$P^o(s'|s) \stackrel{\text{def}}{=} \mathbb{E}_{D:s \rightarrow s'|o} \left[\gamma_d^D \mathbb{1}_{S_D=s'} \right]$$

γ_p controls how much we care about option duration (pseudo-primitive when $\gamma_p=1$)

Options with Time Dilation



Reward model:

$$R^o(s) \stackrel{\text{def}}{=} \mathbb{E}_{D:s|o} \left[\sum_{t=0}^{D-1} \gamma^t R_{t+1} \right]$$



$$R^o(s) \stackrel{\text{def}}{=} \mathbb{E}_{D:s|o} \left[\sum_{t=0}^{D-1} \gamma_r^t R_{t+1} \right]$$

Transition model:

$$P^o(s'|s) \stackrel{\text{def}}{=} \mathbb{E}_{D:s \rightarrow s'|o} \left[\gamma^D [S_D = s'] \right]$$



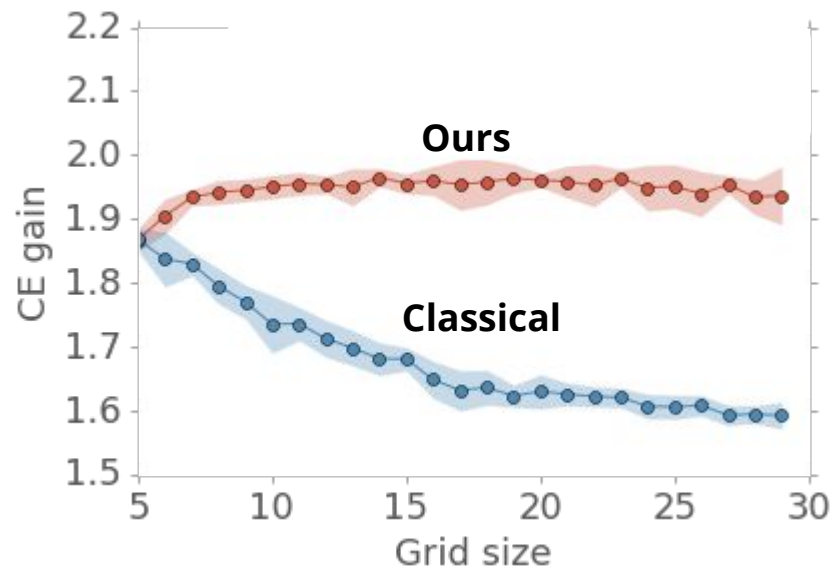
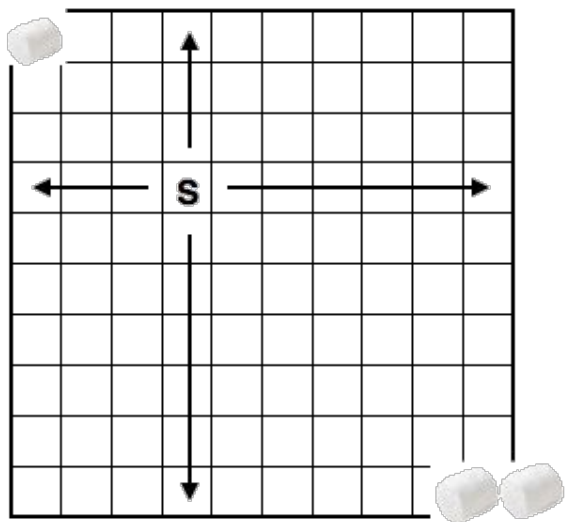
(2) per-decision (1) decouple

$$P^o(s'|s) \stackrel{\text{def}}{=} \gamma_d \mathbb{E}_{D:s \rightarrow s'|o} \left[\gamma_p^D [S_D = s'] \right]$$

γ_p controls how much we care about option duration
(pseudo-primitive when $\gamma_p = 1$)

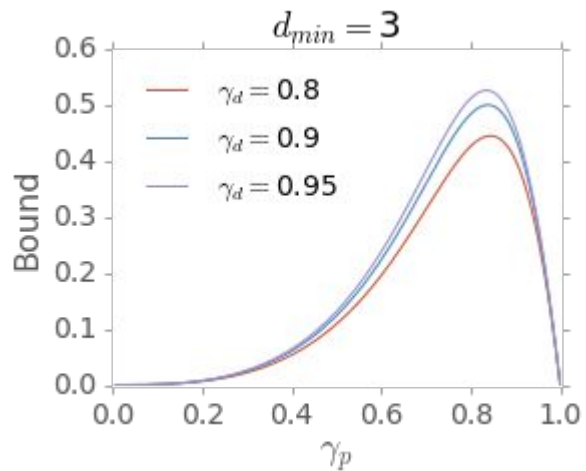
Key intuition: Insulate option time from global time

Primitive Timestep Invariance

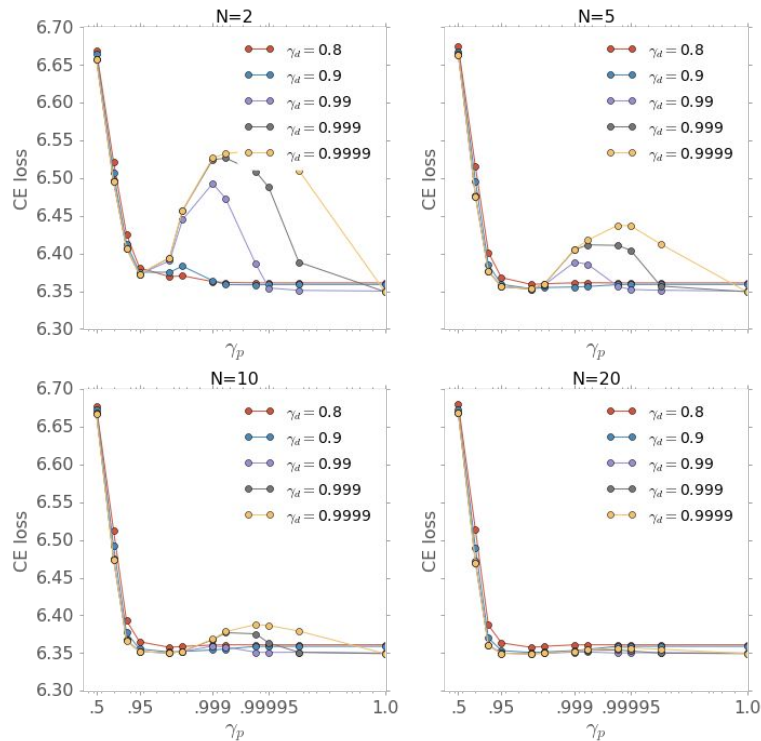


Bias-Variance Tradeoff

Analytical variance bound

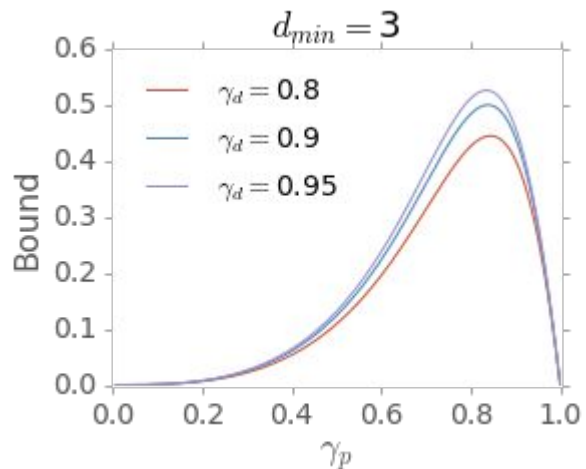


Empirical error (Four Rooms)



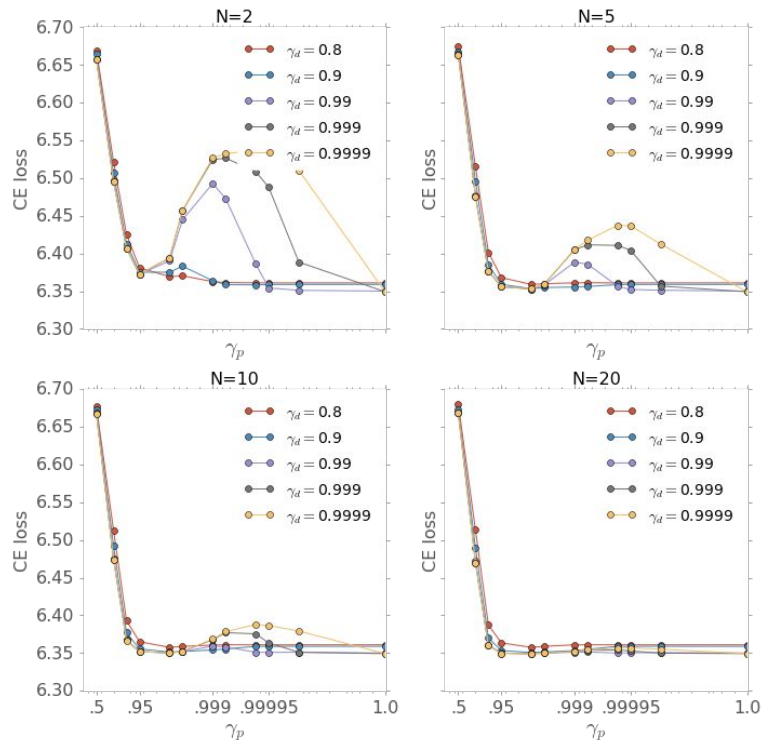
Bias-Variance Tradeoff

Analytical variance bound



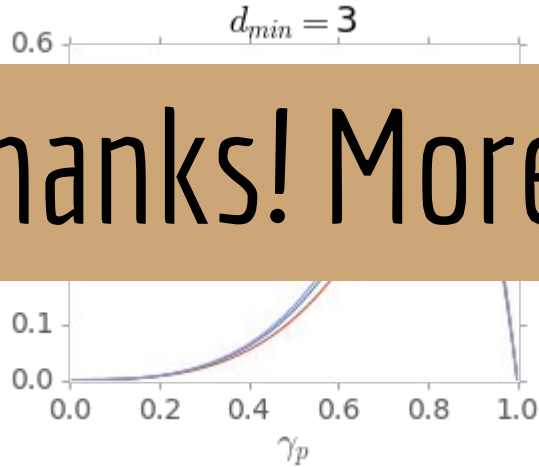
Larger γ_p can induce *less* variance!

Empirical error (Four Rooms)



Bias-Variance Tradeoff

Analytical variance bound



Thanks! More at poster #114 :)

Larger γ_p can induce *less* variance!

Empirical error (Four Rooms)

