5a)

True - gradient TD Learning Rule

$$TD(0): \quad W_{t+1} = W_t + \alpha\left[R_{t+1} + \gamma\hat{v}(s_t', \omega) - \hat{v}(s_t, \omega)\right]\nabla_\omega\hat{v}(s,\omega)$$

on differentiating for $\hat{v}(s', \omega)$ w.r.t $\omega$

$$W_{t+1} = W_t + \alpha\left[R_{t+1} + \gamma\hat{v}(s_t', \omega) - \hat{v}(s_t, \omega)\right]$$

$$\nabla_\omega\hat{v}(s_t, \omega) - \alpha\gamma\nabla_\omega\hat{v}(s_t', \omega)$$

The term $-\alpha\gamma\nabla_\omega\hat{v}(s_t', \omega)$, accounts for the change in the estimated value of next state as $\omega$ changes. which makes this method a true gradient method. As it considers the changes in $\omega$ that affect the feature value estimate.

5b) The learning rule optimizes an objective func'n that includes an expectation over the next state $s'$, making it more similar with mean squared error of the value fun'n estimate., so theoritically

$$\Rightarrow E_\pi\left[(R + \gamma\hat{v}(s', \omega) - \hat{v}(s, \omega))^2\right]$$

would be the objective function.

Incorporating the true could lead to more stable and accurate updates because it accurately represents the objective of minimizing the difference b/w estimated value and true value of a state across all transitions

5c) Mean Squared Bellaman Error (MSBE) objective

For the MBSE objective:

$$BE(\omega) = \sum_{s \in S} \mu(s) \left[ E_\pi [R + \gamma \hat{v}(s', \omega) | S] - \hat{v}(s, \omega) \right]^2$$

To optimize this we differentiate $BE(\omega)$ w.r.t $\omega$ to find the gradient and use it in a gradient descent learning rule:

$$\therefore \nabla_\omega BE(\omega) = -2 \sum_{s \in S} \mu(s) \left( E_\pi [R + \gamma \hat{v}(s', \omega) | S] - \hat{v}(s, \omega) \right) \nabla_\omega \hat{v}(s, \omega)$$

The TD-learning rule that optimizes this objective fun would adjust the weights in the direction that minimizes the $BE(\omega)$, so:

$$\omega_{t+1} = \omega_t - \alpha \nabla_\omega BE(\omega)$$

5d)   2n Code Snippets