# Final Exam Question 4
## STAT 560 Statistical Theory I

Clayton Allard

December 22, 2022

(a) We iteratively add one variable at a time to see which variable makes the biggest improvement based on the t-test. After we add the variable, we run the whole model to see if every variable still remains significant at the $\alpha$ level from the t-test. If one of the variable's P-values fails to be within the threshold, then we take the previous step's model as our choice. If we get all the way through with all variables being significant, then the full model is the best model.

---

**Algorithm 1:** Forward Stepwise Algorithm

---

**Data:** For $n > p$ we have $x \in \mathbb{R}^{n \times p}, y \in \mathbb{R}^n, \alpha \in [0,1]$

$I \leftarrow \mathbf{0} \in \mathbb{R}^p$ ;                    /* $I_i = 1$ if $x_i$ is in the model, 0 otherwise. */

$\mathcal{M} \leftarrow \{1\}$ ;                                    /* Intercept model. */

**while** $I \neq \mathbf{1} \in \mathbb{R}^p$ **do**

   $P \leftarrow \mathbf{1} \in \mathbb{R}^p$ ;                            /* Initialize p-values. */

   **for** $i \in \{1, 2, \ldots, p\}$ **do**

      **if** $I_i = 1$ **then**

         /* $x_i$ is already in the model.                        */

         **continue**

      **end**

      $\mathcal{N} \leftarrow \mathcal{M}$ append $x_i$ ;                    /* Add variable to model. */

      $P_i \leftarrow$ **P-value**$(\mathcal{N}(x_i))$ ;                    /* P-value of t-test for $x_i$. */

   **end**

   $m \leftarrow \arg\min_i P$ ;                                    /* Best variable. */

   $p^* \leftarrow P_m$ ;                                    /* $p^*$ is the lowest P-value. */

   **if** $p^* < \alpha$ **then**

      $\mathcal{M} \leftarrow \mathcal{M}$ append $x_m$;                    /* Add variable to real model. */

      $I_m \leftarrow 1$;                                    /* Include $x_m$ in the model. */

      $V \leftarrow$ **P-value**$(\mathcal{M})$;            /* P-values where $V_i = 1$ if $x_i$ not in model. */

      /* Element wise product $\Rightarrow$ Max of variables in the model.        */

      **if** $\max\{I \odot V\} \geq \alpha$ **then**

         /* We want to ensure that all variables maintain a low P-value.    */

         $\mathcal{M} \leftarrow \mathcal{M}$ remove $x_m$;                /* Then we go with previous model. */

         **return** $\mathcal{M}$

      **end**

   **end**

   **else**

     **return** $\mathcal{M}$

   **end**

**end**

---

(b) To do the backward stepwise algorithm, start with the full model. If every variable's t-test P-value is less than $\alpha$, then that is our model of choice. Otherwise, remove the variable with the

highest P-value. Then we run the model of the remaining predictors. If all the variables come out with a P-value less than $\alpha$ from the t-test, then that is the model we choose. Otherwise, continue this process of removing the variable with the highest P-value. This process may continue until we are all the way down to the null model.