

**Note:** All “@@@\_complete” fields are largely made up of 2’s, with some 1’s and 0’s. Not noted in the data dictionary.

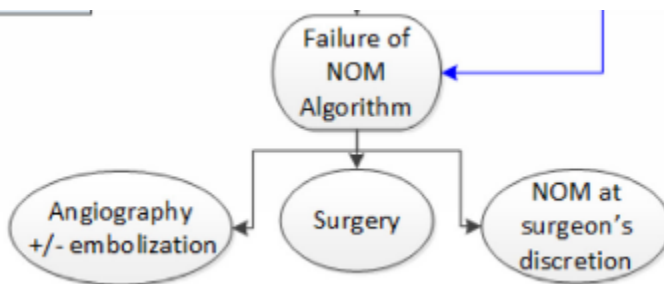
**Note:** Log skips over included features that didn’t need notes.

**Note:** Maybe if less than 5% of non-null data per column, don’t use? (check all text-columns, “other” columns)

## CS180: Manual Feature Selection List

**Removed:** id, redcap\_data\_access\_group

**Include: \*\*\*time\_to\_angio\*\*\*:** might be really significant but not noted in the data dictionary, indicates “failure of NOM Algo” aka an outcome to refer to, lots of NULL - *must ask if NULL means no intervention or just no angio.* **Note:** there is an “angiography” and “embolization” feature



**Removed:** mo\_injury, yr\_injury, ins

**Unsure of significance:** demographics\_complete

**Include: hx\_trauma and hxtrauma\_days:** can be condensed into 1 column since hxtrauma\_days has a range of 1-364. can just use 0 for the other values, hxtrauma\_days only has about 4 non-null numbers out of the 1000 patients. so we should just use hx\_trauma as a true/false feature between known and unknown trauma within the last year OR between definitively no trauma and no known trauma.

**Include: osh\_used:** only a couple unknowns (888), can exclude those patients?, then only use 0-No 1-Yes.

**Unsure if we should use: time\_ed, time\_ptc** (timely response rate might matter)

**Include: moi:** method of injury, one-hot encoding needed for 11 categories +1 unknown (888), but decent amount of “others” noted in text in **moi\_other**

**Include: bleed:** leave as is, only has 0/1 values for no/yes.

*Maybe remove:* **admit** (lots of non-null data, but not sure if necessary)

*Unsure of significance:* **injury\_characteristics\_and\_admission\_complete**

*Include:* **unstable\_elem\_#**: #1-#9 and 0 for None, lots of data noted in text in **other\_elem** like moi

*Unsure if we should use:* **ems\_used** (lots of non-null data, but not sure if necessary)

*Include:* **sbp\_field, dbp\_field, sbp\_lowed, dbp\_lowed, sbp\_lowptc, dbp\_lowptc, hr\_lowfield, hr\_highfield, hr\_lowed, hr\_highed, hr\_lowptc, hr\_highptc, hypo\_ed\*\*, hypo\_aftered\*\*, vaso\_first48, vasos\_# (1-7 w/ 888), vaso\_other\*\*\***

- Deal with many NULL values (can try doing average of columns?)
- ed = emergency dept., ptc = pediatric trauma center
- \*\* maybe need to combine/manipulate these two, mostly 0's but there is a decent amount of data
- \*\*\* can likely get rid of vaso\_other because there's only ~10 non-blank

*Unsure if we should use:* **vaso\_hrs** (little data, might not be enough to use)

*Include:* **intubated\_field, intubated\_ed, which\_ed\_\_# (1,2,888)**

*Unsure of significance:* **vitals\_complete**

*Include:* **add\_injuries\_\_# (1-10), gcs, iss, ais\_head, ais\_chest, ais\_face, ais\_abdomen, ais\_extremity, ais\_external**

- Deal with many NULL values (unsure what method to use)

*Unsure of significance:* **scores\_complete**

*Include:* **hb# and time#**: 1-37, will standardize into 1 column for each interval of hours

*Unsure of significance:* **labs\_complete**

*Include:* **transfusion, mtp, transfusion\_hrs\*\*, transfusion\_local\*\*\***

- \*\*only used if transfusion == 1 (unsure of how to handle remaining nulls)
- \*\*\*only used if transfusion == 1 and is categorical, may need to separate

*Unsure of significance:* **prbc\_24osh, prbc\_osh, prbc\_24ptc, prbc\_ptc, ffp\_24osh, ffp\_osh, ffp\_24ptc, ffp\_ptc, platelets\_24osh, platelets\_osh, platelets\_24ptc, platelets\_ptc, blood\_products\_complete**

- Because of LOTS of NULL data

*Include:* **fast and fast\_results:** results (1-3 w/ 888) rely on fast (0/1 values w/ 888). Maybe combine into 1 column

*Include:* **spleen\_grade, spleen\_fluid, spleen\_ctblush, liver\_grade, liver\_fluid, liver\_ctblush, kidney\_grade\*, pancreas\_grade\***

- Data is NULL if injury is only in spleen/only in liver, not NULL for both if both injured
- \*kidney and pancreas grade largely 0's but there is some data

*Remove:* **embo\_reason** (largely NULL, only has values if there was **embolization**, but even then there's so little data)

*Remove:* \*\*Kelly's thought: the rest of embo: **embo\_type, embo\_other, embo\_tech, embo\_techother**, is not currently important enough to include. If we already know they did a non-NOM on the patient, it's all we really need to know for now. If we have time to include/manipulate, maybe include.

- But they separated these features into multiple columns for us already?
- Largely 0's anyway

*Unsure of significance:* **radiology\_complete**

*Maybe include (likely not):* **dpl, exlap, exlap\_hrs, exlap\_damage, lap, lap\_hrs, lap\_open, exlap\_lap\_purpose**

- Largely 0, expands into **exlap\_hrs, exlap\_damage** with is largely NULL

*Include:* **procedure\_# (0-12), procedure\_other\***

- already separated into columns for each category for us
- \*other category is in text, (unsure of how to handle, maybe remove)

*Maybe remove:* **percu\_drain, percu\_hrs, percu\_fluid\_# (0-7, 888), percu\_other, percu\_performed, percu\_performed\_other, percu\_us**

- Largely 0 or NULL. ~11 data points, maybe not enough data to be significant, even though the feature description seems useful

*Include:* **other\_procedure**

*Remove:* **other\_procedure\_describe**

- Might be too hard to use the "describe" text. But there is decent amount of **other\_procedure** data

*Unsure of significance:* **procedures\_complete**

*Include:* **missed\_iai** (seems like a significant feature through description, even if little data)

*Maybe include:* **describe\_iai** (could be hard to use, and there's little data)

*Remove:* **missed\_ii, describe\_ii** (basically no data, only like 3 datapoints)

*Include:* **missed\_any** (seems like a significant feature through description, even if little data)

*Remove:* **describe\_any** (could be hard to use, and there's little data)

*Remove:* **missed\_treatment** (not very relevant/not much data)

*Include:* **describe\_treatment** (very little data, but keywords of "surgery")

*Include:* **hosp\_comps and comps\_#**: comps (1-14) is separated into columns already, hosp\_comps can be converted to a "none" column for comps (invert 0's and 1's)

*Maybe remove:* **ards\_vent** (not a lot of data), **reor\_operation** (not a lot of data), **trans\_reaction** (all NULL)

*Include:* **comps\_other** (lots of data, but how to handle all text?)

*Include:* **infection\_# (0-6, 888)**, **infection\_other** (lots of data, but how to handle all text?)

*Include:* **disposition** (3 categories, might separate into 3 columns, unless not a lot of "3's")

*Include:* **cause\_death** (mainly 3's and 6's, maybe can make use of), **cause\_other** (mostly bleeding or brain injuries, can manually manipulate)

*Unsure of significance:* **outcomes\_complete**

*Remove:* **pt\_status** (redundant)

*Include:* ALL **stable\_** and **unstable\_**, if NULL, then change to 0 because it's mutually exclusive

*Remove:* **stable\_deviation\_decscribe** (could be unnecessary and very very little data)

*Include:* **unstable\_notprotocol\_other** (a lot of "death" text, may need to manually comb through for consistency)

*Unsure of significance:* **algorithm\_adherence\_complete**

*Remove:* **fu1\_intervention\_procedure, fu1\_radiology, fu1\_radabd, fu1\_abdet, fu1\_abdet\_finding\_# (1-4), fu1\_abdet\_other** (largely NULL or 0's)

**Decision:** disregard all follow-ups (**fu1\_** and **fu2\_**). Doesn't really solve the problem at hand, whether the patient is actively in shock and needs surgery. fu#\_intervention and fu#\_intervention\_procedure not actually necessary.

*Include:* **hgt, bmi, firsttemp\_ptc, lowtemp\_ptc**, (~50% NULLs and 888s, but must deal with)

*Remove:* **firsttemp\_ptc\_method, lowtemp\_ptc\_method, lasttemp\_ptc, lasttemp\_ptc\_method, firsttemp\_icu, firsttemp\_icu\_method, time\_normtemp**

- Just so much missing data, mostly normal temperatures/similar temperatures to firsttemp\_ptc and lowtemp\_ptc

*Unsure of significance:* **additions\_complete**

**Note:** *Rest of columns in the dataset seems to be used for their own KNN prediction. May not be relevant to ours, or not enough information to be using those columns.*