

# Penalized logistic regression on time-to-event data using casebase sampling

---

Jesse Islam

2/27/2020

# Popular methods in time-to-event analysis

- In disease etiology, we tend to make use of the proportional hazards hypothesis.
  - Cox Regression
- When we want the absolute risk:
  - Breslow estimator
  - Parametric models

## Motivations for a new method

- Julien and Hanley found that survival analysis rarely produces prognostic functions, even though the software is widely available in cox regression packages. [2]
- They believe the stepwise nature is the reason, as it reduces interpretability. [2]
- A streamlined approach for reaching a **smooth absolute risk** curve. [2]

## Dr. Cox's perspective

**Reid:** How do you feel about the cottage industry that's grown up around it [the Cox model]?

**Cox:** Don't know, really. In the light of some of the further results one knows since, I think I would normally want to tackle problems parametrically, so I would take the underlying hazard to be a Weibull or something. I'm not keen on nonparametric formulations usually.

**Reid:** So if you had a set of censored survival data today, you might rather fit a parametric model, even though there was a feeling among the medical statisticians that that wasn't quite right.

**Cox:** That's right, but since then various people have shown that the answers are very insensitive to the parametric formulation of the underlying distribution [see, e.g., Cox and Oakes, Analysis of Survival Data, Chapter 8.5]. And if you want to do things like predict the outcome for a particular patient, it's much more convenient to do that parametrically.

1. SUPPORT study
2. Casebase sampling
3. Penalized logistic regression on survival data
4. Maximum likelihood with regularization
5. Absolute risk comparison
6. Conclusion
7. Future work
8. References

- **Study to Understand Prognoses and Preferences for Outcomes and Risks Treatments**
- Design: Prospective cohort study.
- Setting: 5 academic care centers in the United States.
- Participants: 9105 hospitalized.
- Follow-up-time: 5.56 years.
- 68% incidence rate.

## SUPPORT manual imputation [4]

- Notorious for missing data

Baseline Variable	Normal Fill-in Value
Bilirubin	1.01
BUN	6.51
Creatinine	1.01
PaO2/FiO2 ratio (pafi)	333.3
Serum albumin	3.5
Urine output	2502
White blood count	9 (thousands)

**Table 1:** Suggested imputation values. [3]

- Mice imputation package (R)
  1. PMM (Predictive Mean Matching) – For numeric variables
  2. logreg(Logistic Regression) – Binary Variables
  3. polyreg(Bayesian polytomous regression) Factor Variables



## Removed variables [4]

- **Response variables**
  - Hospital Charges.
  - Patient ratio of costs to charges.
  - Patient Micro-costs.
- **Ordinal covariates**
  - functional disability.
  - Income.
- **Sparse covariate**
  - Surrogate activities of daily living.
- Previous model results. (6)

- **Response variables**
  - follow-up time, death.
- **Covariates**
  - Age, sex, race, education (6)
  - Disease group/class, comorbidities. (3)
  - Coma score, Therapeutic Intervention Scoring System (2)
  - Physiological variables. (11)
  - Activities of daily living. (2)

## Original SUPPORT analysis [4]

- Determined SUPPORT prognostic model on phase I patients.
- Tested on Phase II.
- Both on the scale of 180 days.

## Original SUPPORT analysis [4]

SUPPORT physiology score (SPS) was developed.

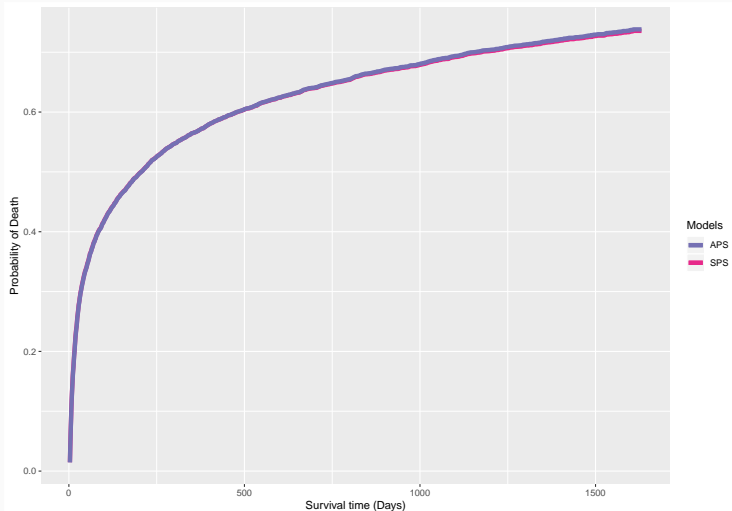
$$\begin{aligned} \text{SPS} = & 259.9\{\text{ARF/MOSF}\} + 263.4\{\text{COPD/CHF}\} + \\ & 241.4\{\text{Cirrhosis/Coma}\} + 281.5\{\text{Lung/Colon Cancer}\} - \\ & 0.06174 \min(\text{PaO}_2/\text{FiO}_2, 225) - 0.6316 \min(\text{Mean BP}, 60) \\ & + 1.0205 \text{ WBC} - 0.3676(\text{WBC} - 8)_+ - 0.5631(\text{WBC} - \\ & 11)_+ + 0.2691 \min(\text{Alb}, 4.6) + 0.2312 \text{ Aresp} - 2.362 \\ & \text{Temp} + 1.326(\text{Temp} - 36.6)_+ + 2.473(\text{Temp} - 38.3)_+ \\ & - 1.579 \times 10^{-1} \text{ HR} + 9.770 \times 10^{-5} (\text{HR} - 55)_+^3 - 2.189 \\ & \times 10^{-4} (\text{HR} - 80)_+^3 + 1.518 \times 10^{-4} (\text{HR} - 110)_+^3 - \\ & 3.062 \times 10^{-5} (\text{HR} - 149)_+^3 + 0.9763 \text{ Bil} - 0.7481(\text{Bil} - \\ & 7)_+ - 6.8761 \text{ Cr} + 11.6058(\text{Cr} - 0.600)_+^3 - 21.8413(\text{Cr} - \\ & 1.000)_+^3 + 10.3574(\text{Cr} - 1.500)_+^3 - 0.1219(\text{Cr} - \\ & 5.399)_+^3 - 0.6167096 \text{ Na} + 0.0021118(\text{Na} - 128)_+^3 - \\ & 0.0036730(\text{Na} - 135)_+^3 + 0.0006126(\text{Na} - 139)_+^3 + \\ & 0.0009486(\text{Na} - 148)_+^3 - 6.278 \{\text{COPD/CHF}\} \times \min \\ & (\text{Alb}, 4.6) - 11.45 \{\text{Lung/Colon Cancer}\} \times \min(\text{Alb}, \\ & 4.6) + \{\text{ARF/MOSF}\}[-2.3549 \text{ WBC} + 2.7494 (\text{WBC} - \\ & 8)_+ - 0.4638 (\text{WBC} - 11)_+] \end{aligned}$$

- How does their SPS perform over 5.56 years?
- How does the Apache III physiology score (APS) perform over 5.56 years?
- How does a model with all the covariates, excluding SPS and APS, perform?

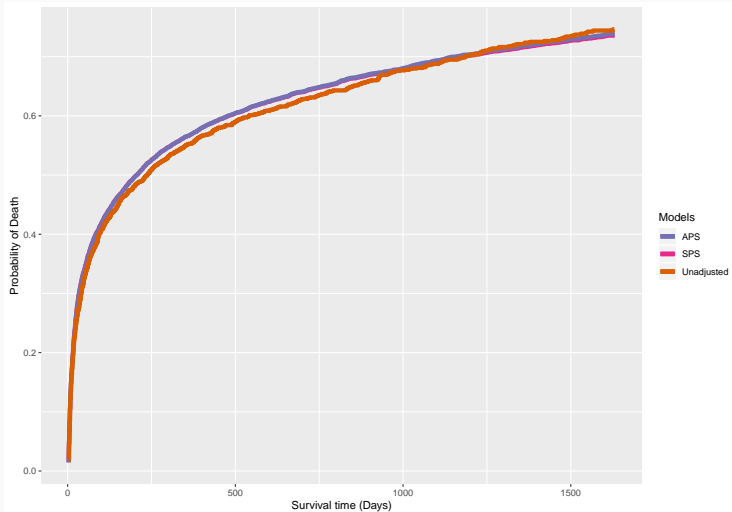
# Analysis Process

1. Impute
  2. Compare SPS and APS over ~5.56 years using absolute risk.
  3. Compare to Kaplan-Meier curve (unadjusted model).
  4. Compare to full model (excluding SPS and APS).
- All models is trained on 80% of the observations.
  - Remaining observations are used to generate comparative absolute risk curves.
    - The absolute risk curve for each individual is averaged.
    - Averaged curve is expected to approach Kaplan-meier curve.

# SPS vs APS

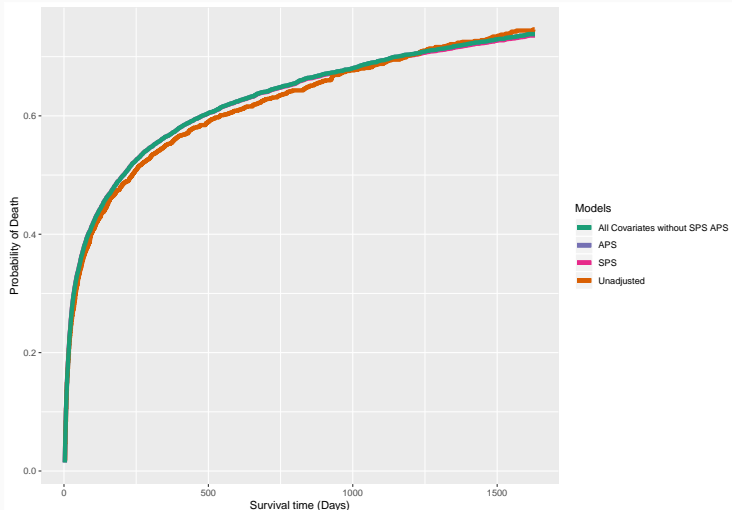


# SPS vs. Kaplan-Meier

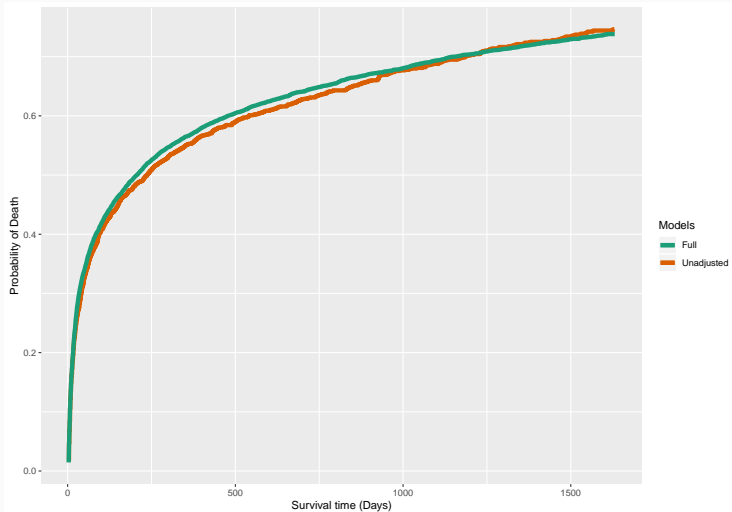




# All covariates vs. physiology scores vs unadjusted



# Chosen absolute risk comparisons



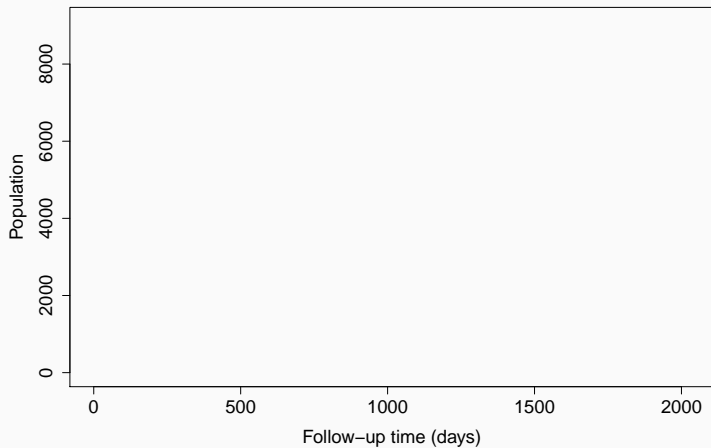
## Chosen absolute risk comparisons: conclusion

- Linear associations without physiology scores perform similarly to SPS and APS alone.
- We choose the linear associations without physiology scores as the model of choice (Full model).

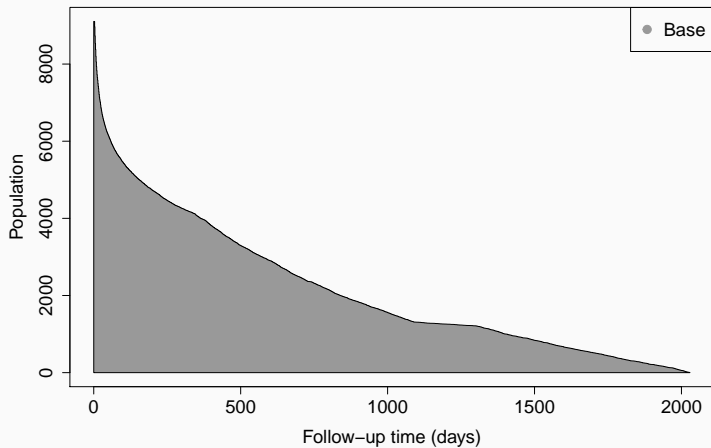
# Casebase sampling overview

1. Clever sampling.
  2. Implicitly deals with censoring.
  3. Allows a parametric fit using *logistic regression*.
- Casebase is parametric, and allows different parametric fits by incorporation of the time component.
  - Package contains an implementation for generating *population-time* plots.

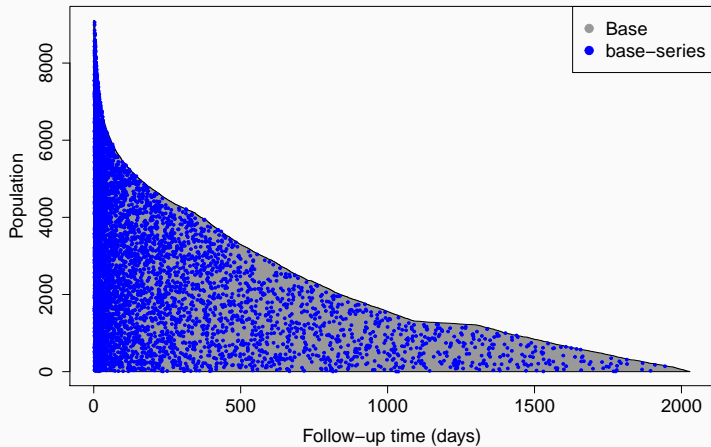
## Casebase: Sampling



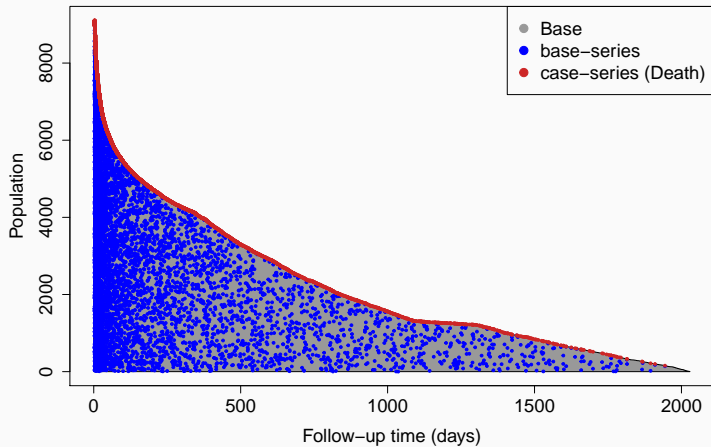
## Casebase: Sampling



## Casebase: Sampling

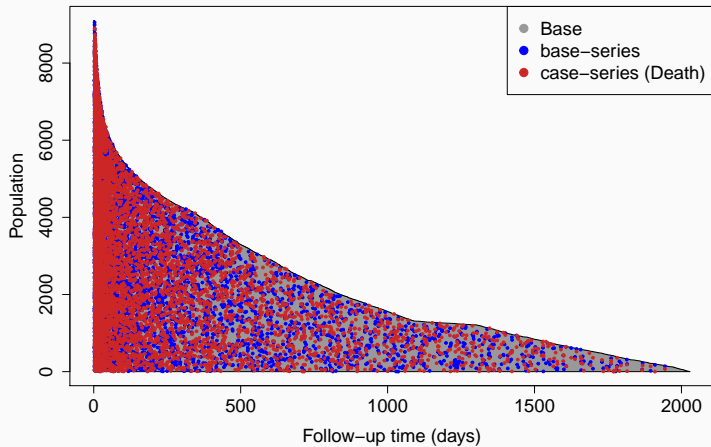


## Casebase: Sampling





# Casebase: Sampling

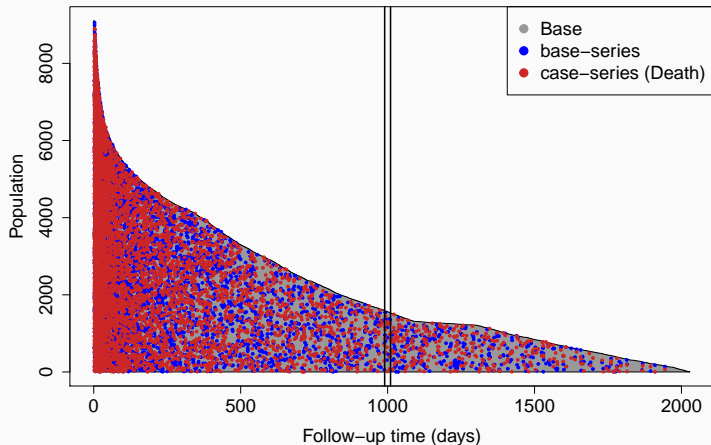


$$e^L = \frac{Pr(Y = 1|x, t)}{Pr(Y = 0|x, t)} = \frac{h(x, t) * B(x, t)}{b[B(x, t)/B]} = \frac{h(x, t) * B}{b}$$

- $L = \beta X$
- $b$  = base-series.
- $B$  = Base.
- $B(x, t)$  = Risk-set for survival time  $t$ .

# Casebase: Sampling

$$e^L = \frac{Pr(Y=1|x,t)}{Pr(Y=0|x,t)} = \frac{h(x,t)*B(x,t)}{b[B(x,t)/B]} = \frac{h(x,t)*B}{b}$$



## log-odds = log hazard

$$e^L = \frac{\hat{h}(x, t) * B}{b}$$

$$\hat{h}(x, t) = \frac{b * e^L}{B}$$

$$\log(\hat{h}(x, t)) = L + \log\left(\frac{b}{B}\right)$$

## Maximum log-likelihood [1]

$$\log(l(\beta_0, \beta)) = \frac{1}{N} \sum_{i=1}^N \{y_i(\beta_0 + x_i^T \beta) - \log(1 + e^{\beta_0 + x_i^T \beta})\}$$

## Maximum log-likelihood, with offset

$$\log(l(\log(\frac{b}{B}), \beta)) = \frac{1}{N} \sum_{i=1}^N \{y_i(\log(\frac{b}{B}) + \mathbf{x}_i^T \beta) - \log(1 + e^{\log(\frac{b}{B}) + \mathbf{x}_i^T \beta})\}$$

## Maximum log-likelihood, with offset and lasso

$$\frac{1}{N} \sum_{i=1}^N \{y_i(\log(\frac{b}{B}) + x_i^T \beta) - \log(1 + e^{\log(\frac{b}{B}) + x_i^T \beta})\} - \lambda \|\beta\|$$

- We can now fit models of the form:

$$\log(h(t; \alpha, \beta)) = g(t; \alpha) + \beta X$$

- By changing the function  $g(t; \alpha)$ , we can model different parametric families easily:



## Casebase: Parametric models

*Exponential:*  $g(t; \alpha)$  is equal to a constant

```
casebase::fitSmoothHazard(status ~ X1 + X2)
```

*Gompertz:*  $g(t; \alpha) = \alpha t$

```
casebase::fitSmoothHazard(status ~ time + X1 + X2)
```

*Weibull:*  $g(t; \alpha) = \alpha \log(t)$

```
casebase::fitSmoothHazard(status ~ log(time) + X1 + X2)
```

- We have a bunch of different parametric hazard models now.
- To get the absolute risk, we need to evaluate the following equation in relation to the hazard:

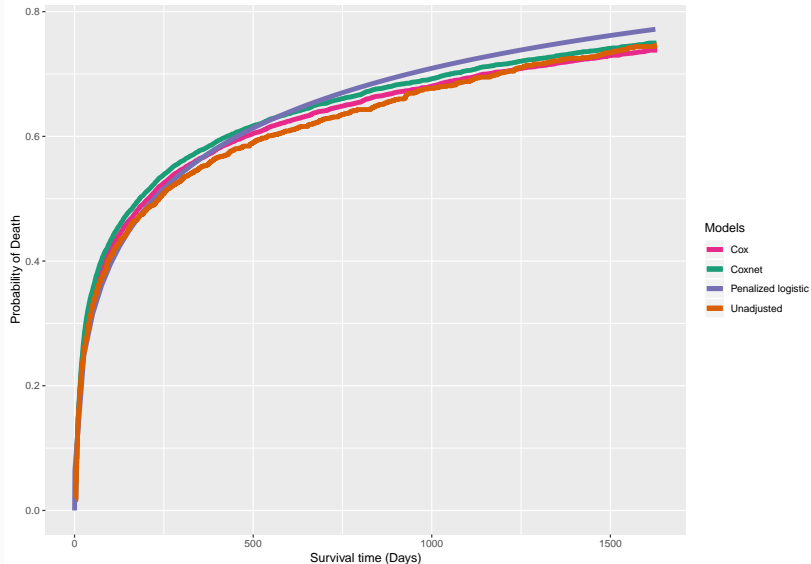
$$CI(x, t) = 1 - e^{-\int_0^t h(x, u) du}$$

- $CI(x, t)$  = Cumulative Incidence (Absolute Risk)
- $h(x, u)$  = Hazard function
- Lets use the weibull hazard.

## Models to be compared

- Recall: Case study to demonstrate regularization using our method.
- **unadjusted**:  $(\text{death}, \text{time}) \sim 1$
- **Cox**:  $(\text{death}, \text{time}) \sim \beta X$
- **Coxnet**:  $(\text{death}, \text{time}) \sim \beta X \leftarrow \text{Lasso}$
- **Penalized logistic**:  $\text{death} \sim \log(\text{time}) + \beta X \leftarrow \text{Lasso}$

# Survival comparison



- Classical survival analysis requires methods to incorporate censorship in our data.
- Case-base sampling is a technique that implicitly incorporates censorship implicitly.
- Logistic regression on SUPPORT dataset had slightly different results near the end of follow-up time.

- Comparative measure.
- Survival GWAS.

## References 1

- 1.Czepiel, S. A. (2002). Maximum likelihood estimation of logistic regression models: theory and implementation. Available at [czep.net/stat/mlelr.pdf](http://czep.net/stat/mlelr.pdf), 1825252548-1564645290.
- 2.Hanley, James A, and Olli S Miettinen. 2009. "Fitting Smooth-in-Time Prognostic Risk Functions via Logistic Regression." *The International Journal of Biostatistics* 5 (1).
- 3.Harrell, F. (2020). SupportDesc < Main < Vanderbilt Biostatistics Wiki. [online] [Biostat.mc.vanderbilt.edu](http://biostat.mc.vanderbilt.edu/wiki/Main/SupportDesc). Available at: <http://biostat.mc.vanderbilt.edu/wiki/Main/SupportDesc> [Accessed 25 Feb. 2020].

## References 2

- 4.Knaus, W. A., Harrell, F. E., Lynn, J., Goldman, L., Phillips, R. S., Connors, A. F., . . . & Layde, P. (1995). The SUPPORT prognostic model: Objective estimates of survival for seriously ill hospitalized adults. *Annals of internal medicine*, 122(3), 191-203.
- 5.Saarela, Olli, and Elja Arjas. 2015. "Non-Parametric Bayesian Hazard Regression for Chronic Disease Risk Assessment." *Scandinavian Journal of Statistics* 42 (2). Wiley Online Library: 609–26.
- 6.Saarela, Olli. 2015. "A Case-Base Sampling Method for Estimating Recurrent Event Intensities." *Lifetime Data Analysis*. Springer, 1–17



7. Scrucca L, Santucci A, Aversa F. Competing risk analysis using R: an easy guide for clinicians. *Bone Marrow Transplant*. 2007 Aug;40(4):381-7. doi: 10.1038/sj.bmt.1705727.
8. Turgeon, M. (2017, June 10). Retrieved May 05, 2019, from <https://www.maxturgeon.ca/slides/MTurgeon-2017-Student-Conference.pdf>

## **Tutorial:**

<http://sahirbhatnagar.com/casebase/>

## **Slides:**

[https://github.com/Jesse-Islam/ATGC\\_survival\\_presentation\\_Feb.27.2020](https://github.com/Jesse-Islam/ATGC_survival_presentation_Feb.27.2020)

## **Questions?**