



# Topics in Machine Learning

## Machine Learning for Healthcare

Rahul G. Krishnan  
Assistant Professor

Computer science & Laboratory Medicine and Pathobiology

# Last week

---

- Supervised machine learning
- Risk stratification
  - Stratification as a prediction problem
  - **Case study:** Predicting the onset of diabetes



# Outline

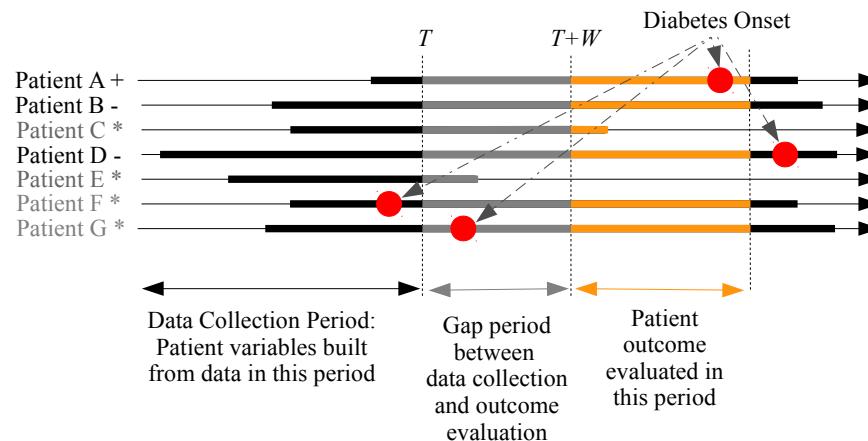
- Clarification to questions
- Risk stratification
  - Deriving labels
  - Evaluating models
- Survival analysis:
  - From binary to continuous valued outcomes
  - Parametric
  - Non-parametric
  - Semi-parametric

# Questions from last week

- Is there structure among diagnosis codes:
  - Yes! They are organized in a hierarchy. View ICD10 codes [here](#)

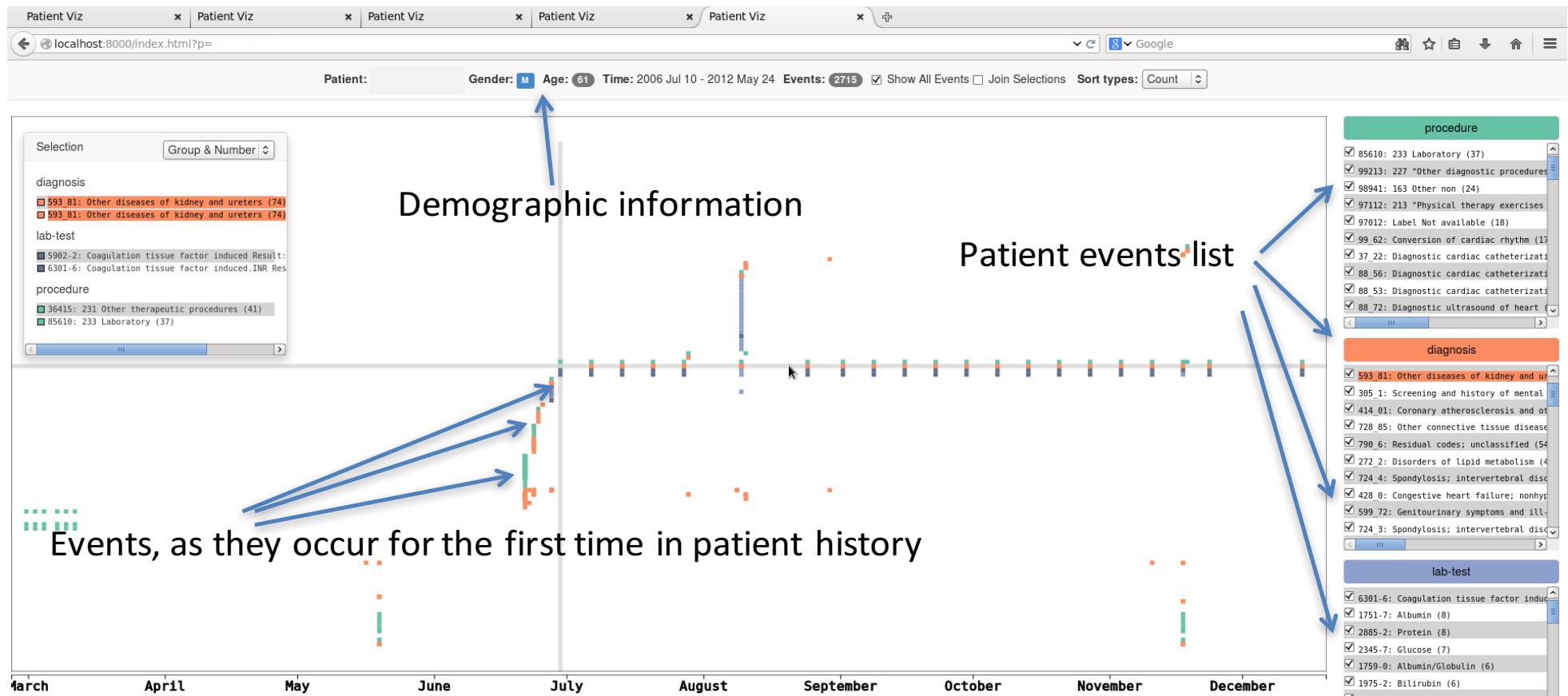
[ICD-10-CM Codes](#) > [E00-E89 Endocrine, nutritional and metabolic diseases](#) > [E08-E13 Diabetes mellitus](#) >  
Type 2 diabetes mellitus E11  
**Type 2 diabetes mellitus E11-**
- Using predictive models in the future (on different data):
  - Non-stationarity of data is a challenging problem
    - Means that data distribution changes in unpredictable ways over time
  - Covariate shift can tank good machine learning models deployed in clinics
  - Still a lot of research on good techniques for detection of covariate shift

# Deriving labels for risk stratification



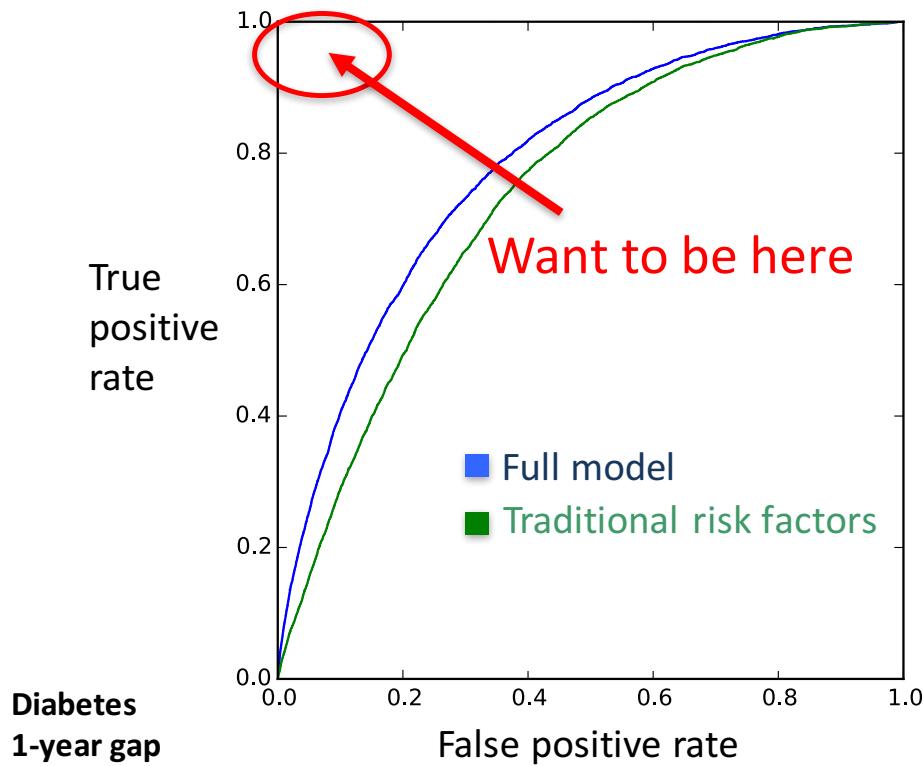
- Typically done via chart review
- Work with doctor to assess criteria that constitute **Diabetic Onset**
  - e.g. does the patient have ICD10 code for diabetes

# Evaluation of risk stratification models



<https://github.com/nyuviz/patient-viz>

# Evaluation of risk stratification models



**AUC = Area under the ROC curve**

- Invariant to class imbalance
- Interpretable as the probability that an algorithm ranks a positive patient over a negative patient

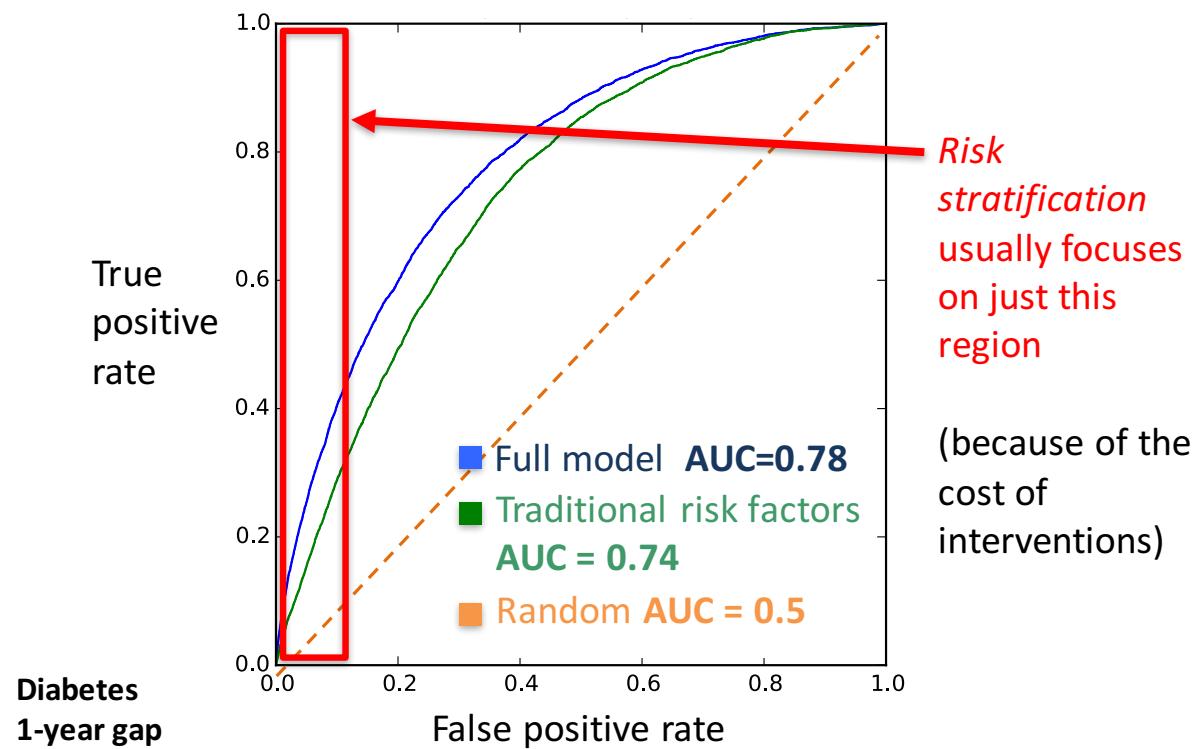
$$TPR = \frac{TP}{Actual\ Positive} = \frac{TP}{TP + FN}$$

$$FNR = \frac{FN}{Actual\ Positive} = \frac{FN}{TP + FN}$$

$$TNR = \frac{TN}{Actual\ Negative} = \frac{TN}{TN + FP}$$

$$FPR = \frac{FP}{Actual\ Negative} = \frac{FP}{TN + FP}$$

# Where you want to be on ROC curve



# Many other important statistical considerations when building models

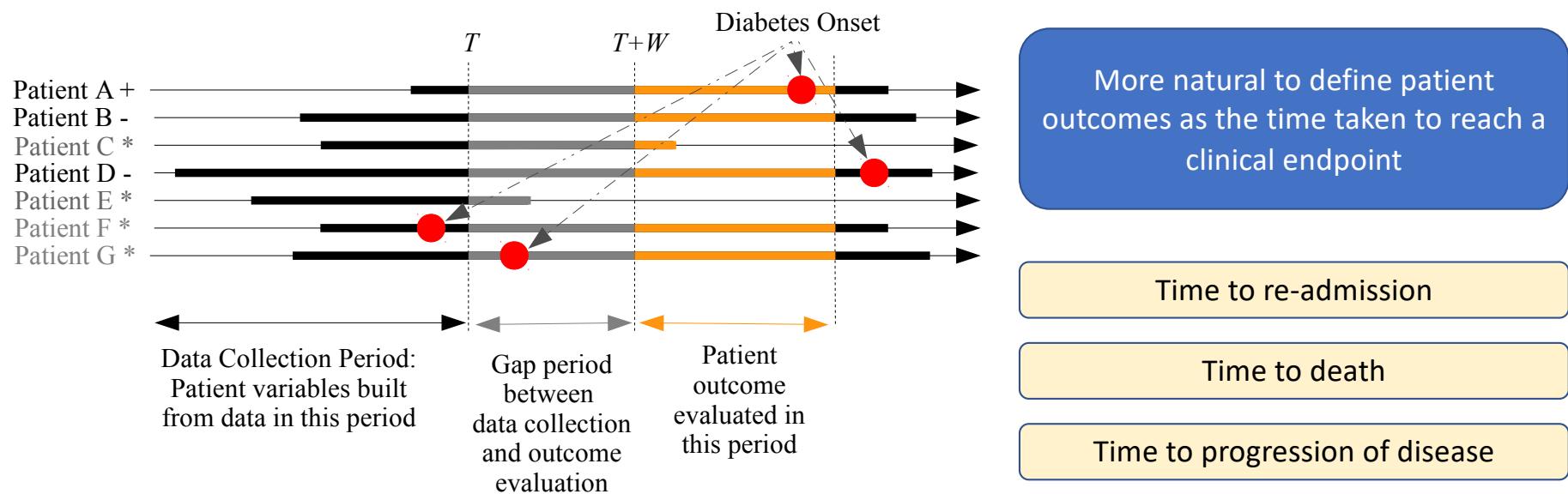
- Calibration
- Sensitivity analysis
- Error bars and confidence intervals on prediction estimates
- Heterogeneity of results:
  - Does the model only work well for a subpopulation?
- Model introspection:
  - For a linear model, are the features used by the models the ones you might expect?
  - Do root nodes in a decision tree make sense?
  - More challenging to do for deep neural networks

# The importance of interpretability

- [Intelligible Models for HealthCare: Predicting Pneumonia Risk and Hospital 30-day Readmission](#), Caruana et. Al, KDD 2015
- Used generalized additive models to make predictions of pneumonia and readmission
- Learn  $\text{HasAsthma}(x) = \text{LowerRiskOfDying}(x)$
- Why?
  - Asthmatics w/ pneumonia are prioritized
  - Get aggressive treatment, faster, in ICU
  - Treatment lowers risk of death compared to general population
- Scenario where the prescription of an intervention taints the outcomes
- The consequence:
  - Automated methods might flag asthmatics as not being problematic!

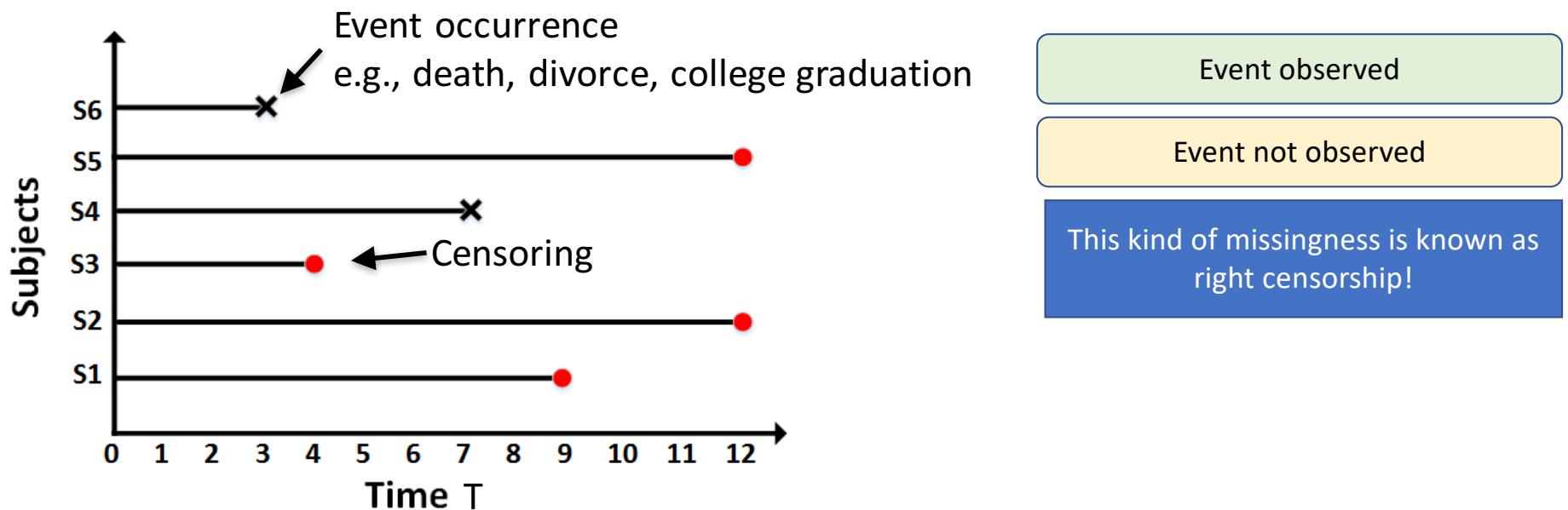
Questions?

# Continuous valued outcomes for risk stratification



# Labels are partially missing!

- We may not observe the outcome in the dataset – life, administrative challenges etc.



# Censorship

- Censorship is an important to know about when handling longitudinal data
- Three types of censorship:
  - Left censored data:
    - We don't observe the *start* of an event but we do observe longitudinal data after it
    - Example: ICU patient's vitals are continuously recorded from when they enter. If one of the sensors fails and is later fixed, their data is left censored
  - Right censored data [focus for today]
    - We don't observe the *incidence* of an event but we know it occurs after the last observed time
    - We'll see examples soon!
  - Interval censored data:
    - Both left and right censorship
    - Neonatal unit is tracking data on children, observe data sometime after they are born and until they leave the unit (for those who survive)

# What can we do when we do not observe when the event occurs?

- What do we know:
  - $x$ : features
  - $y$ : last observed time
  - $b$ : whether or not the event occurs
- Option: why not throw away all datapoints for which we don't know when the event occurs:
  - Wasteful, might end up with very little data
- Key idea behind survival analysis:
  - Learn to predict time-to-event with all the available data that we have

# Preliminaries

- $(x, T, b) = (\text{features}, \text{time}, \text{censoring})$ 
  - $b = 0$  if censored and  $b = 1$  if event is observed
- $f(t) = p(t) = \text{probability of death at time } t$
- Survival function:  $S(t) = P(T > t) = \int_t^\infty f(x)dx$
- Hazard function:  $h(t) = \lim_{\epsilon \rightarrow 0} p(T \in (t, t + \epsilon] | T \geq t)$
- Cumulative hazard function:  $H(t) = \int_0^t h(u)du = \int_0^t \frac{f(u)}{S(u)}du = \int_0^t \frac{-dS(u)}{S(u)}du = -\log\{S(t)\}$
- Hazard function & survival function:  $h(t) = \frac{f(t)}{S(t)}$ 
$$f(t) = h(t)S(t) = h(t) \exp\{-H(t)\}$$
$$S(t) = \exp\{-H(t)\}.$$

Slide credit: Lu Tian and Richard Olshen's course on survival analysis

# Preliminaries – (2)

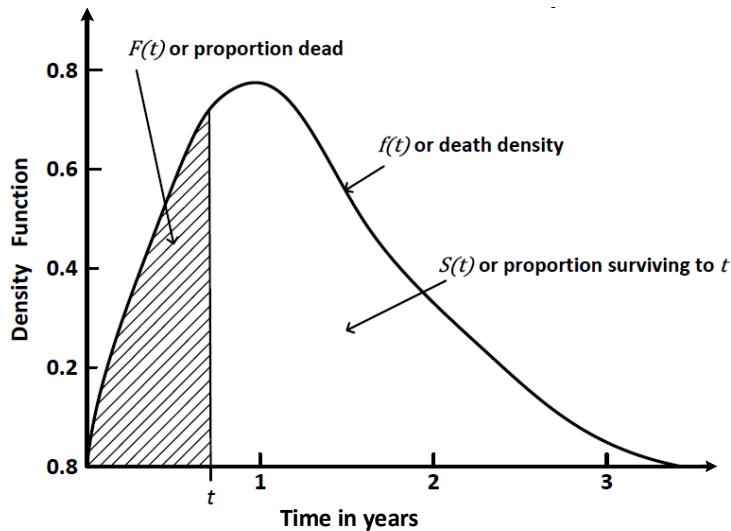


Fig. 2: Relationship among different entities  $f(t)$ ,  $F(t)$  and  $S(t)$ .

[Wang, Li, Reddy. Machine Learning for Survival Analysis: A Survey. 2017]

[Ha, Jeong, Lee. Statistical Modeling of Survival Data with Random Effects. Springer 2017]

# Preliminaries

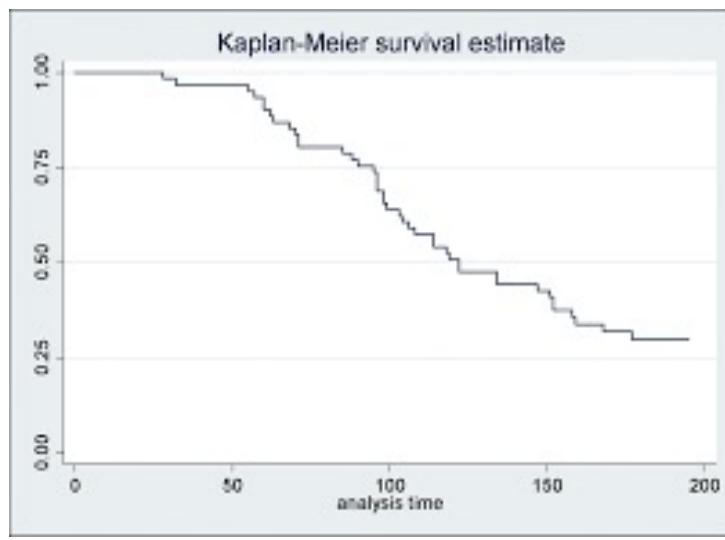
- $h(t) dt$  is approximately the conditional probability of the event occurring in an infinitesimal interval around  $t$  conditional on it not having occurred before  $t$

# Non-parametric survival analysis

- Let start by ignoring our features and asking about  $S(t)$ 
  - $S(t)$  is an integral 
$$S(t) = P(T > t) = \int_t^{\infty} f(x)dx$$
  - **Idea:** If we had access to  $f(x)$  we could discretize time and evaluate  $f(x)$  in each bin and sum them up.
  - **Issue:** We don't have access to  $f(x)$  but we **do** have samples!
- **Kaplan Meier curves:**
  - **Non-parametric** estimator of the survival function  $S(t)$ 
    - We do not assume anything about the underlying distribution of  $S(t)$
    - We'll use our entire dataset to approximate the shape of  $S(t)$

# Kaplan Meier estimator

- Derivation out of scope for this class
  - Survival analysis is a rich area of research and is often a course in and of itself
  - E.g. [Lu Tian and Richard Olshen at Stanford](#)



Observed event times

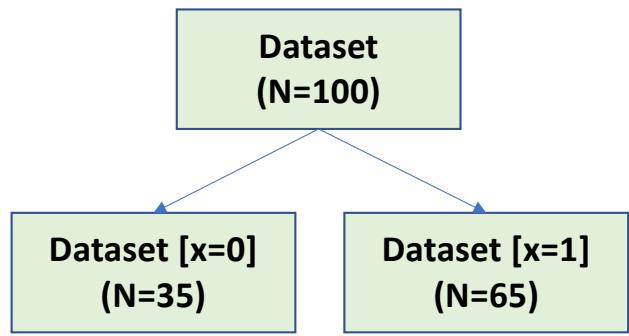
$$y_{(1)} < y_{(2)} < \dots < y_{(D)}$$

$d_{(k)}$  = # events at this time

$n_{(k)}$  = # of individuals alive  
and uncensored

$$\widehat{S}_{K-M}(t) = \prod_{k:y_{(k)} \leq t} \left\{ 1 - \frac{d_{(k)}}{n_{(k)}} \right\}$$

# What do we do if we have features ( $x$ )?



Evaluate KM estimator on each strata

Right [survival probability of patients who have multiple myeloma stratified by genetic marker]

Survival probability,  
 $S(t)$

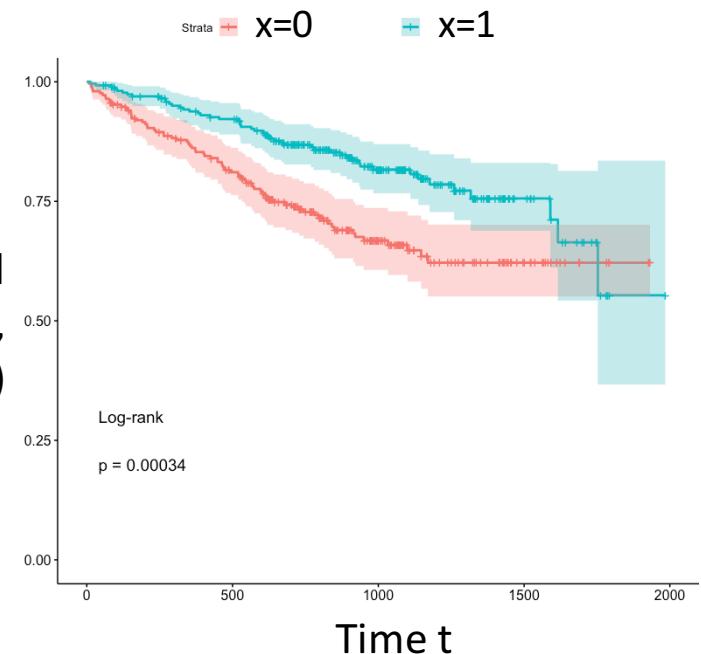


Figure credit: Rebecca Boiarsky

# What if $x$ is high-dimensional?

- Option 1: stratify based on clusters of  $x$
- Option 2: Let the survival function **depend on  $x$**
- Lets study what happens in supervised learning:
- In linear regression:  $y \sim \mathcal{N}(w^T x + b; 1)$ 
  - Outcome is a Gaussian function centered around  $w^T x + b$
  - Known as a parametric model for  $y$ :
    - There are some parameters that govern the behavior of  $y$  as a function of  $x$

# Maximum likelihood estimation for supervised learning

$x_1$   $y_1$

$x_2$   $y_2$

$x_3$   $y_3$

Dataset (N=3)

- Given a dataset, the model parameters are learned via **maximum likelihood estimation**

$$\mathcal{L}(y, x) = \log p(y|x; \theta)$$

Score function (high is good, low is bad)

$$\theta = \arg \max_{\theta} \sum_{i=1}^N \mathcal{L}(y_i, x_i)$$

Solve this optimization problem to **learn** the model. Often formulated as a minimization of the negative of the log-likelihood function

# Maximum likelihood estimation for survival analysis

$x_1$   $t_1$   $b_1$

$x_2$   $t_2$   $b_2$

$x_3$   $t_3$   $b_3$

Dataset (N=3)

- Given a dataset, the model parameters are learned via **maximum likelihood estimation**

$$p(T = t|x; \theta) = f(t) \quad \text{Uncensored likelihood}$$

$$p(T > t|x; \theta) = S(t) \quad \text{Censored likelihood}$$

$$\sum_{i=1}^N b_i \log p(T = t_i|x_i; \theta) + (1 - b_i) \log p(T > t_i|x_i; \theta)$$

Maximize the following objective function to learn model parameters

# What distribution should I use?

**Table 2.1** Useful parametric distributions for survival analysis

Distribution	Survival function $S(t)$	Density function $f(t)$
Exponential ( $\lambda > 0$ )	$\exp(-\lambda t)$	$\lambda \exp(-\lambda t)$
Weibull ( $\lambda, \phi > 0$ )	$\exp(-\lambda t^\phi)$	$\lambda \phi t^{\phi-1} \exp(-\lambda t^\phi)$
Log-normal ( $\sigma > 0, \mu \in R$ )	$1 - \Phi\{(\ln t - \mu)/\sigma\}$	$\varphi\{(\ln t - \mu)/\sigma\}(\sigma t)^{-1}$
Log-logistic ( $\lambda > 0, \phi > 0$ )	$1/(1 + \lambda t^\phi)$	$(\lambda \phi t^{\phi-1})/(1 + \lambda t^\phi)^2$
Gamma ( $\lambda, \phi > 0$ )	$1 - I(\lambda t, \phi)$	$\{\lambda^\phi / \Gamma(\phi)\}t^{\phi-1} \exp(-\lambda t)$
Gompertz ( $\lambda, \phi > 0$ )	$\exp\{\frac{\lambda}{\phi}(1 - e^{\phi t})\}$	$\lambda e^{\phi t} \exp\{\frac{\lambda}{\phi}(1 - e^{\phi t})\}$

(parameters  
can be a  
function of  $x$ )

[Ha, Jeong, Lee. Statistical Modeling of Survival Data with Random Effects. Springer 2017]

# CoxPH: Interpretability in survival analysis

- Parameteric models that depend on  $x$  change parameters of a distribution in linear/non-linear ways as a function of  $x$
- **Goal:** Link variation to covariates directly to the survival function
- The [Cox Proportional Hazard's model](#) is one of the most popular tools in survival analysis

$$h(t|X = x; \theta) = \underbrace{h_0(t)}_{\text{Baseline hazard}} \exp(\beta^T x)$$

Baseline hazard reflects the hazard for subjects with all covariates equal to 0

## Interpretation in the univariate case

$$\frac{h(t|X = x_1; \theta)}{h(t|X = x_2; \theta)} = \frac{\exp(\beta^T x_1)}{\exp(\beta^T x_2)} \quad \frac{h(t|X = x + 1)}{h(t|X = x)} = \exp(\beta)$$

Hazard ratio is independent of time

Parameters have an intuitive meaning

CoxPH: Linear model for log of the hazard ratio

# CoxPH for binary data

- $X = [\text{received drug (0 no, 1 yes)}, \text{gender (0 male, 1 female)}]$

$$h(t|x_1, x_2) = h_0(t) \exp(\beta_1 z_1 + \beta_2 z_2)$$

$$h(t|X) = h_0(t)$$

No treatment

Male

$$h(t|X) = h_0(t) \exp(\beta_1)$$

Yes treatment

Male

$$h(t|X) = h_0(t) \exp(\beta_2)$$

No treatment

Female

$$h(t|X) = h_0(t) \exp(\beta_1 + \beta_2)$$

Yes treatment

Female

# Key advantage of the CoxPH model

- We can estimate the model parameters  $\beta$
- Notably we can do so without estimation the baseline hazard
- This is a **semi-parametric model**
  - We make no assumptions about the baseline hazard rate but we do learn parameters corresponding to how it is modified based on patient covariates
- How do we learn this model?
  - Won't derive from scratch
  - Useful reference: [Course notes by Ronghui \(Lily\) Xu](#)

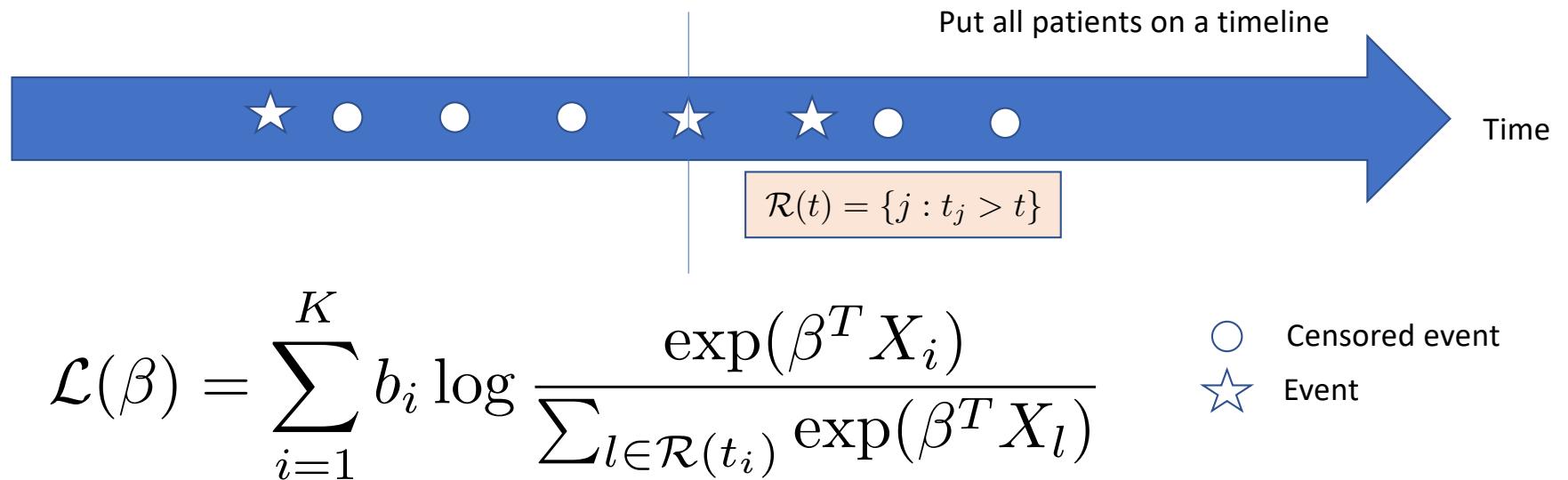
# Cox Partial Likelihood

Scan from left to right in time, at each discrete point, calculate the risk set

$$\mathcal{R}(t) = \{j : t_j > t\}$$

$$\mathcal{L}(\beta) = \sum_{i=1}^K b_i \log \frac{\exp(\beta^T X_i)}{\sum_{l \in \mathcal{R}(t_i)} \exp(\beta^T X_l)}$$

# Visualizing the computation of the partial likelihood



Intuition: How likely are the features of this patient to explain their elevated risk of having the event occur now compared to all the individuals whose event occurs later!

# Advances in machine learning for survival analysis

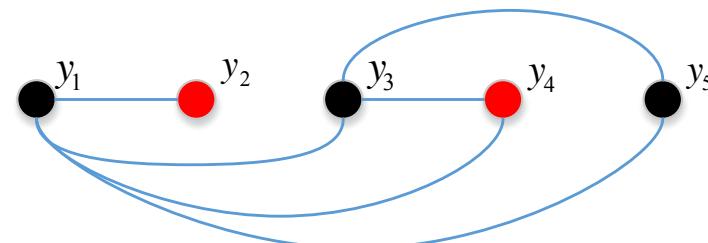
- [DeepSurv, Katzman et. al, 2017](#)
  - One of the readings for this week uses a deep neural network to parameterize the modification to the hazard function
  - Parameter estimation by taking derivatives of the
- Advanced reading [Deep survival analysis, Ranganath et. al, 2016](#)
  - What if  $x$  is very high dimensional?
  - Rather than condition on  $x$  directly, learn a latent representation of  $x$  while jointly modeling survival time

# Evaluation in survival analysis

- Concordance index (aka C-statistic) – predicts how well the model ranks patients based on survival (i.e. predicts relative survival time)
- Equivalent to AUC (when there is no censoring)

$$\hat{c} = \frac{1}{num} \sum_{i:b_i = 0} \sum_{j:y_i < y_j} I[S(\hat{y}_j|X_j) > S(\hat{y}_i|X_i)]$$

Black = uncensored  
Red = censored



## Other ways to evaluate models

- Mean squared error [for just those who are uncensored]
- Held out likelihood (censored + uncensored)

# Questions?

---



## Why not use classification?

If T has a non-uniform distribution,  
Thresholds for classification may not be known  
at training time



## Why not use regression?

When outcomes are missing [event time not observed] you may have to throw data out

- Leads to limited training data
- Might introduce bias into the dataset

# Announcements

- Piazza now has a pinned post to help you start looking for teammates to work on problems
- Watch MIMIC data and analysis tutorials [posted to course website]
- Friday:
  - 30 minutes - Discussion with [Alistair Johnson, Scientist, SickKids Hospital](#)
  - 20 minutes - Breakout rooms to brainstorm ideas for group projects