# Nearest Neighbor and Kernel Survival Analysis

Nonasymptotic Error Bounds and Strong Consistency Rates

George H. Chen Assistant Professor of Information Systems Carnegie Mellon University



June 11, 2019



### **Survival Analysis**



**Goal:** Estimate  $S(t|x) = \mathbb{P}($ survive beyond time  $t \mid$  feature vector x)

## **Problem Setup**

**Model:** Generate data point  $(X, Y, \delta)$  as follows:



# Theory (Informal)

*k*-NN estimator with  $k = \widetilde{\Theta}(n^{2\alpha/(2\alpha+d)})$  has strong consistency rate:

$$\sup_{t \in [0,\tau]} |\widehat{S}(t|x) - S(t|x)| \le \widetilde{O}(n^{-\alpha/(2\alpha+d)})$$

If no censoring, problem reduces to conditional CDF estimation

→ Error upper bound, up to a log factor, matches conditional CDF estimation lower bound by Chagny & Roche 2014

Proof ideas also give finite sample rates for:

- Kernel Kaplan-Meier estimators
- k-NN & kernel Nelson-Aalen cumulative hazard estimators  $(-\log S(t \mid x))$
- Generalization bound for automatic *k* using validation data

Most general finite sample theory for *k*-NN and kernel survival estimators Existing kernel results only for Euclidean space (Dabrowska 1989, Van Keilegom & Veraverbeke 1996, Van Keilegom 1998)

#### Experiments

