# Part 6: Duration Analysis

Chris Conlon

Microeconometrics

June 15, 2017

# Overview

Another way to look at dynamic behavior

- In the previous lecture we looked a single agent dynamics where optimizing agents made decisions over multiple periods.
- The maintained assumption in that case was that $F(X'|X)$ was a 1st order Markov Process and was exogenous to the agent's decision and we treated it like a nuisance parameter.
- Sometimes the object of interest is actually the transition function itself.
- We are often interested in the length of spells or how much time is spent in each state which we call duration.

# Overview

Simple cases:

- ▶ The simplest cases are single irreversible transitions
    - ▶ Alive → Dead
    - ▶ Working → Failure
- ▶ Other easy cases are "resetting" processes:
    - ▶ Employed → unemployed for zero weeks, one week, etc.
    - ▶ Healthy → Sick Day 1, Sick Day 2, etc.
    - ▶ Not on strike → Strike Day 1, Strike Day 2, etc.
- ▶ Let's start with these before we worry about multivariate outcomes or more complicated cases.

# Decisions

Have to make some decisions first

1. Do we model spell length directly or probability of transition?
   - Most of the time we want to work with probability of transition.
   - If we work with probability of transition, we have to pay attention to frequency
2. What outcomes do we measure: stocks? or flows?
   - Do we measure the number of people who lose/find jobs?
   - Do we measure the number of unemployed people each month?
3. Is the data truncated or censored?
   - People who are still alive are not in the dataset!

For now we will think about single-spells, and measure them using flow data.

# Examples

There are lots of different names (depending on your discipline):

- ▶ Life table analysis
- ▶ Hazard Analysis
- ▶ transition analysis
- ▶ survival analysis
- ▶ failure time analysis

Examples:

- ▶ How long does a government last?
- ▶ How long does a part last?
- ▶ How long before a firm adopts a new technology?
- ▶ How long do marriages last?
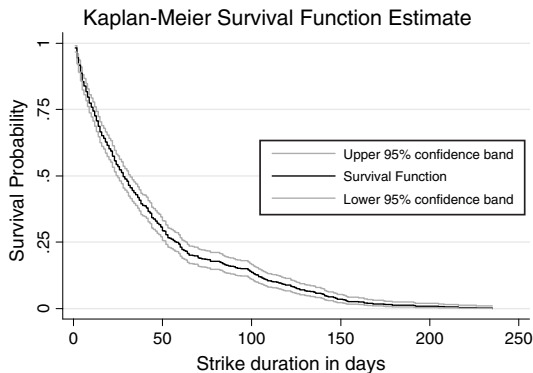- ▶ How long before criminals re-offend?

# Start with a Graph!



**Figure 17.1:** Strike duration: Kaplan-Meier estimate of survival function. Data on completed spells for 566 strikes in the U.S. during 1968–76.

# What did we just plot?

The empirical survival function

- We ignored any covariates, including calendar time.
- The x-axis was the duration
- The y-axis was the fraction of observations still alive "alive" after $x$ periods.
- If nothing is infinitely lived then the graph always starts at $1$ and always ends at zero.
- If things are infinitely lived we call the duration distribution defective.

# Parametric

Let's start with some deeply parametric stuff

- density function: $f(t) = dF(t)/dt$: unconditional probability of instantaneous failure
- CDF: $F(t) = Pr(T \leq t) = \int_0^\infty f(s)ds$. (Probability that spell is less than length $t$).
- Survival Function: $S(t) = 1 - F(t) = Pr(T > t)$. This has the nice property that it integrates to expected duration $\int_0^\infty S(t)dt = E[T]$.
- Hazard Function: $\lambda(t) = \lim_{\Delta t \to 0} \frac{Pr[t \leq T < t + \Delta t | T \geq t]}{\Delta t} = \frac{f(t)}{S(t)}$.
- All of these functions represent the same information!

# More about hazard functions

- Hazard is conditional probability of leaving unemployment after being unemployed for $t$.
- Hazard is percentage change in survivor function $S(t)$
- Hazard also gives us the distribution of duration $T$:

$$\lambda(t) = -\frac{\partial \log S(t)}{\partial t}$$

$$S(t) = \exp\left[-\int_0^\infty \lambda(u)du\right]$$

- Often we'd like to estimate $\lambda(t|x)$ instead of $E[T|x]$ especially since we often have censored data so that $\lambda(t|x)$ is still well defined but $E[T|x]$ is not.
- In practice $\lambda(t|x)$ can be tricky to estimate (especially since it may contain zeros at some $t$ in finite sample. Solution: Cumulative Hazard Function.

$$\Lambda(t) = \int_0^\infty \lambda(s)ds = -\log S(t)$$

- Just like we preferred to estimate CDF instead of PDF!

# Summary Table

**Table 17.1.** *Survival Analysis: Definitions of Key Concepts*

| Function | Symbol | Definition | Relationships |
|----------|--------|------------|---------------|
| Density | $f(t)$ | | $f(t) = dF(t)/dt$ |
| Distribution | $F(t)$ | $\Pr[T \leq t]$ | $F(t) = \int_0^t f(s)ds$ |
| Survivor | $S(t)$ | $\Pr[T > t]$ | $S(t) = 1 - F(t)$ |
| Hazard | $\lambda(t)$ | $\lim_{h \to 0} \dfrac{\Pr[t \leq T < t+h \mid T \geq t]}{h}$ | $\lambda(t) = f(t)/S(t)$ |
| Cumulative hazard | $\Lambda(t)$ | $\int_0^t \lambda(s)ds$ | $\Lambda(t) = -\ln S(t)$ |

# What about Discrete Time?

- Maybe we only see survival annually/weekly/etc. not actual failure time.
- Basic idea is the same. Have to be careful about ties. Divide failures into $t_j$ buckets

$$
\begin{aligned}
\lambda_j &= Pr[T = t_j | T \geq t_j] = f^d(t_j)/S^d(t_{j-}) \\
\Lambda^d(t) &= \sum_{j | t_j \leq t} \lambda_j \\
S^d &= Pr[T \geq t] = \prod_{j | t_j \leq t} (1 - \lambda_j)
\end{aligned}
$$

- Can define the product integral which is regular product in discrete case and exponential of integral in continuous case.

# Nonparametric estimation

- Without censoring, things are easy: just let

$$\hat{S}(t) = \frac{1}{n} \sum_{i=1}^{n} \mathbf{1}(T_i \geq t).$$

- if you want a smooth hazard function, take a smooth estimator, e.g. (with some "small" bandwidth $w > 0$)

$$\hat{S}(t) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{1 + \exp((t - T_i)/w)},$$

- and then take minus the derivative of the log of this estimate.

What if there is censoring? Kaplan-Meier!

# Censoring

Lots of ways for censoring to arise:

- ▶ Mostly concerned about right censoring
- ▶ Observe spells from time $0$ until time $c$ and all we know is that they end in $(c, \infty)$.

# Kaplan-Meier

- We define the ordered durations as

$$T_{(1)} < \ldots < T_{(n)},$$

- let $d_j$ be the number of observations $i$ for which $T_i = T_{(j)}$
- Let $m_j$ number of spels censored in $[t_j, t_{j+1})$
- and $r_j$ the cardinality of the risk set at duration $t_{j-}$
  $r_j = \sum_{l|l \geq j} d_l + m_l$
- Simple estimate of the hazard function $\hat{\lambda}_j = \frac{d_j}{r_j}$.
- Kaplan-Meier estimator of the survival function is the Product Limit Estimator

$$\hat{S}(t) = \Pi_{j|t_j \leq t} \left( 1 - \frac{d_j}{r_j} \right) = \Pi_{j|t_j \leq t} \left( \frac{r_j - d_j}{r_j} \right)$$

- It is normally distributed (asymptotically), with (Greenwood) variance

$$\hat{V}[\hat{S}(t)] = (\hat{S}(t))^2 \cdot \sum_{j|t_j \leq t} \frac{d_j}{r_j(r_j - d_j)}.$$

# Other stuff

Think about what happens when $m_j = 0$ (no censoring)

- $r_j = \sum_{l|l \geq j} d_l + m_l \rightarrow r_{j+1} = r_j - d_j$.
- $\hat{S}(t) = \Pi_{j|t_j \leq t} \left( \frac{r_j - d_j}{r_j} \right) = \Pi_{j|t_j \leq t} \frac{r_{j+1}}{r_j} = \frac{r_j}{N}$
- Again – exactly what we would expect – one minus the ECDF.

How do we deal with ties?

- Lots of ties can create problems. Implicitly we assume all deaths are at same time in period.
- Why does this matter– well how many are remaining in $r_j$?
- $r_j$ is potentially biased if we have lots of ties.
- Can either try corrections or sample data at higher frequency

# Parametric models

Usually we specify directly the hazard function (closer to theory).

- *Economic example:* an unemployed person ($T \geq t$) leaves unemployment when (given that his benefits decrease in time) he has an offer with a wage $w \geq r(t)$, a reservation wage

- wage offers arrive with some probability $p_t$, and a distribution such that $\Pr(w \geq s) = \bar{G}_t(s)$.

- Then someone unemployed at $t$ leaves unemployment at $t$ with probability $p_t \bar{G}_t(r(t))$ so the hazard function is

$$\lambda(t) = p_t \bar{G}_t(r(t)).$$

- A job search model would give us a theory for $r(.)$, $p_t$ and $\bar{G}_t$, up to some parameters to be estimated, and conditional on covariates $x$.

# Exponential and Weibull

- The exponential is popular because it has a constant hazard rate $\lambda(t) = \gamma$ which does not depend on $t$.

- This is often referred to as the memorylessness property of the exponential.

- This is analytically convenient but it makes it hard to fit things in practice (you only have one parameter!)

- The Weibull is a generalization with $\lambda(t) = \gamma \alpha t^{\alpha-1}$. For $\alpha = 1$ we have exponential.

- For $\alpha > 1$ i is increasing and for $\alpha < 1$ it is decreasing (monotonically).

- Weibull used to be popular in econometrics for simple parametric analysis.

# Exponential and Weibull

**Table 17.4.** *Exponential and Weibull Distributions: pdf, cdf, Survivor Function, Hazard, Cumulative Hazard, Mean, and Variance*

| Function | Exponential | Weibull |
|----------|-------------|---------|
| $f(t)$ | $\gamma \exp(-\gamma t)$ | $\gamma \alpha t^{\alpha-1} \exp(-\gamma t^{\alpha})$ |
| $F(t)$ | $1 - \exp(-\gamma t)$ | $1 - \exp(-\gamma t^{\alpha})$ |
| $S(t)$ | $\exp(-\gamma t)$ | $\exp(-\gamma t^{\alpha})$ |
| $\lambda(t)$ | $\gamma$ | $\gamma \alpha t^{\alpha-1}$ |
| $\Lambda(t)$ | $\gamma t$ | $\gamma t^{\alpha}$ |
| $E[T]$ | $\gamma^{-1}$ | $\gamma^{-1/\alpha}\Gamma(\alpha^{-1}+1)$ |
| $V[T]$ | $\gamma^{-2}$ | $\gamma^{-2/\alpha}[\Gamma(2\alpha^{-1}+1)-[\Gamma(\alpha^{-1}+1)]^2]$ |
| $\gamma, \alpha$ | $\gamma > 0$ | $\gamma > 0, \alpha > 0$ |

# Comparison of Parametric Models

**Table 17.5.** *Standard Parametric Models and Their Hazard and Survivor Functions*[a]

| Parametric Model | Hazard Function | Survivor Function | Type |
|---|---|---|---|
| Exponential | $\gamma$ | $\exp(-\gamma t)$ | PH, AFT |
| Weibull | $\gamma \alpha t^{\alpha-1}$ | $\exp(-\gamma t^\alpha)$ | PH, AFT |
| Generalized Weibull | $\gamma \alpha t^{\alpha-1} S(t)^{-\mu}$ | $[1 - \mu \gamma t^\alpha]^{1/\mu}$ | PH |
| Gompertz | $\gamma \exp(\alpha t)$ | $\exp(-(\gamma/\alpha)(e^{\alpha t} - 1))$ | PH |
| Log-normal | $\dfrac{\exp\left(-(\ln t - \mu)^2/2\sigma^2\right)}{t\sigma\sqrt{2\pi}[1 - \Phi((\ln t - \mu)/\sigma)]}$ | $1 - \Phi\left((\ln t - \mu)/\sigma\right)$ | AFT |
| Log-logistic | $\alpha \gamma^\alpha t^{\alpha-1}/[(1 + (\gamma t)^\alpha)]$ | $1/[1 + (\gamma t)^\alpha]$ | AFT |
| Gamma | $\dfrac{\gamma(\gamma t)^{\alpha-1} \exp[-(\gamma t)]}{\Gamma(\alpha)[1 - I(\alpha, \gamma t)]}$ | $1 - I(\alpha, \gamma t)$ | AFT |

[a] All the parameters are restricted to be positive, except that $-\infty < \alpha < \infty$ for the Gompertz model.

# Adding Covariates

- We can also add covariates by letting $\gamma = \beta X$.
- Sometimes this is called link function or generalized linear models similar to what we saw with the logit or probit.
- It is usually a bad idea to link more than one nonlinear parameter this way.
- We would typically estimate via MLE. Writing down the full-data log-likelihood is straightforward.
- A frequently used special-case are proportional hazard models

# The Proportional Hazard Model

With covariates $x$, the hazard function is $h(t|x)$; we specify

$$\lambda(t|x) = \lambda_0(t)\phi(x).$$

- $\lambda_0$ and $\phi$ are up to a positive multiplicative constant.)
- We call $\lambda_0$ the baseline hazard; every individual has a hazard that is just a proportional version of the baseline hazard.

The baseline hazard could be:

- constant: the survival function is exponential
- a power function $\lambda_0(t) = \gamma t^{\alpha}$; e.g. for $\alpha < 0$ we have negative duration dependence (the long-term unemployed. . . )
- more complicated (flexible) specifications.

# Estimating the PH Model

**Maximum likelihood:** works for any parametric model $\lambda(t|x, \beta)$ of the full hazard function;
(here: w/o censoring, without corrleation across individuals):

$$\max_{\beta} \sum_{i=1}^{n} \ln f(T_i|x_i, \beta),$$

where $f(t|x, \beta)$ is the density of the duration $T$ induced by $\lambda$:

$$f(t|x) = \lambda(t|x)S(t|x) = \lambda(t|x) \exp(-\Lambda(t|x)),$$

so the log-likelihood for $i$ is just $\ln \lambda(T_i|x_i, \beta) - \Lambda(T_i|x_i, \beta)$.

## What's the point?

- ▶ The (partial) additive separability of the log-likelihood in the PH model is designed to make our lives easier.
- ▶ Presumably, we specified $\lambda$ so that its integral $\Lambda$ is easy to compute.
- ▶ For PH: the log-likelihood for $i$ is:
  $\ln \lambda_0(T_i, \beta) + \ln \phi(x_i, \beta) - \Lambda_0(T_i, \beta)\phi(x_i, \beta)$.
- ▶ The most common choice is $\phi(x_i, \beta) = \exp(x_i\beta)$ so that $\ln \phi(x_i, \beta) = x_i\beta$.
- ▶ In that case we have that $\partial\lambda/\partial x_j = \beta_j \cdot \lambda$.
- ▶ One remaing problem: what to do with the baseline hazard function (is that even identified?).

# Cox's Partial Likelihood for the PH Model

- if we do not want to assume anything about the shape of the baseline hazard function
- but we are happy specifying $\phi(x, \beta)$
- then we will only look at the *order* of the durations: we reorder individuals so that $T_{i_1} < \ldots < T_{i_n}$
- ...and we forget about the durations! Then the partial likelihood is:

$$\sum_{j=1}^{n} \left( \ln \phi(x_{i_j}, \beta) - \ln \left( \sum_{l=j}^{n} \phi(x_{i_l}, \beta) \right) \right).$$

- This is a limited information maximum likelihood estimator. It is not fully efficient!
- But it may be robust to mis-specifying $\lambda_0$. Is it actually a valid likelihood? not sure!.

# How did that work?

Once we have ordered everything:

- Let $R(t_j)$ be the set of spells at risk (still alive) at $t_j$
- $d_j$ are the deaths at time $t_j$ $\sum_l \mathbf{1}[t_l = t_j]$.
- Consider only at-risk spells ending a fixed $t_j$

$$
\begin{aligned}
Pr[T_j = t_j | R(t_j)] &= \frac{Pr[T_j = t_j | T_j \geq t_j]}{\sum_{l \in R(t_j)} Pr[T_l = t_l | T_l \geq t_j]} \\
&= \frac{\lambda_j(t_j | x_j, \beta)]}{\sum_{l \in R(t_j)} \lambda_l(t_j, x_l, \beta)} \\
&= \frac{\phi(x_j, \beta)}{\sum_{l \in R(t_j)} \phi(x_l, \beta)}
\end{aligned}
$$

- $\lambda_0$ drops out because of PH.

# Why?

- *Intuition:* those individuals who exit first are (on average) those in the risk set whose covariates $x$ give them the largest $\phi(x, \beta)$.

- After we have $\hat{\beta}$ we can estimate the baseline integrated hazard; denoting $N(t)$=number of individuals with $T = t$

$$\widehat{\Lambda}_0(T_{i_j}) = \sum_{m=1}^{j} \frac{N(T_{i_m})}{\sum_{l=m}^{n} \phi(x_{i_l}, \hat{\beta})}.$$

# Tricks

*A simple way to test the model*:

- just take two different groups of individuals, estimate PH on each, check whether the baseline hazards look <span style="color:red">proportional</span> **NOT** <span style="color:red">equal</span>

*testing a parametric specification of the baseline hazard $\bar{\Lambda}_0$:*

- define generalized residuals $\bar{u}_i = \bar{\Lambda}_0(T_i)$
- Under the true model, for any $z$

$$\Pr(\bar{u} < z) \simeq \Pr(T_i < \bar{\Lambda}_0^{-1}(z)) = 1 - S_0(\bar{\Lambda}_0^{-1}(z)).$$

- it should be $1 - \exp(-z)$ if $S_0 = \exp(-\bar{\Lambda}_0)$.
- So you can estimate the integrated hazard of $(\bar{u}_1, \ldots, \bar{u}_n)$; it should be $\Lambda_u(z) \equiv z$.

# The PH Model is Usually too Restrictive

- ▶ **Fact:** the hazard rate of leaving unemployment decreases in time;
- ▶ It could be *skimming*: the more able, more willing, better connected find a job faster;
- ▶ or it could be "technological": skills deteriorate over time.
- ▶ Under the PH model it can only be the latter: negative duration dependence. $\rightarrow$ introduce unobserved heterogeneity:

$$\lambda(t|x,v) = \lambda_0(t)\phi(x)v.$$

- ▶ $v$ is a "type" that is unobserved by the econometrician; we only assume that it is uncorrelated with $x$ and independent of $t$.

# Dynamic selection

- The model with $v$ is called the **Mixed PH model** (MPH).
- In the unemployment story: the larger $v$'s have a higher hazard rate, so they find a job faster
- Over time, the distribution of $v$ moves (stochastically) to the left.
- This dynamic selection is a general phenomenon in the MPH model: $\lambda(t|x)$ has "more negative duration dependence" than $\lambda(t|x, v)$.
- Can we test dynamic selection vs true negative duration dependence ($\lambda_0$ decreasing)? $\rightarrow$ identification issues.
- This idea shows up in dynamic models of durable goods purchases as well.

## Identification

We still can recover the aggregate survival function from the data, but now it is a mixture:

$$S^A(t|x) = \Pr(T \geq t|x) = \int \exp(-v\phi(x)\Lambda_0(t))dF(v).$$

- Can we recover $\phi$ and $\lambda_0$ without assuming anything on $F$?
- Almost ... in theory: we just need to assume that $E(v)$ is finite.

## A Constructive Proof, 1

- Normalize $Ev = 1$; and $\phi(x_0) = 1$ for some $x_0$.
- Then the aggregate hazard function is

$$\lambda^A(t|x) = -\frac{\partial \log S^A}{\partial t}(t|x)$$

that is

$$\frac{\int v\phi(x)\lambda_0(t)\exp(-v\phi(x)\Lambda_0(t))dF(v)}{S^A(t|x)}.$$

- Look at $x = x_0$ and $t = 0^+$: then $\Lambda_0(t) \simeq 0$, so

$$\lambda^A(0^+|x_0) = \frac{Ev \times k(x_0) \times \lambda_0(0)}{S^A(0|x_0)} = \lambda_0(0).$$

- and

$$\phi(x) = \frac{\lambda^A(0^+|x)}{\lambda^A(0^+|x_0)}.$$

# A Constructive Proof, 2

- Now we can define

$$m^A(t|x) = -\frac{\partial \log S^A(t, x)}{\partial \phi(x)}$$

- and we get the baseline hazard from

$$\frac{\lambda_0(t)}{\Lambda_0(t)} = \frac{\lambda^A(t|x)}{m^A(t|x)};$$

- and we can also recover $F$.
- In practice we would specify functional forms of course.

# Is that Practical?

- We are relying heavily on "identification at 0": that is where we get $\phi(x)$, the rest depends on it.
- Empirical researchers have found that it is often a slim basis (and a very slow-converging estimator)—but anything else will be parametric.
- The alternative is to use richer data: multiple durations/multiple spells.

# Application 1: job search

E.g. Cahuc/Postel-Vinay-Robin, *Econometrica* 2006.

- ▶ Workers are heterogeneous, so are firms;
- ▶ a worker quits when he gets a better outside offer (exogenous Poisson($\lambda$)).
- ▶ We observe (given matched employer-employee data):
  - ▶ job durations (how long each worker stays in a job)
  - ▶ and distributions of wages (mostly) across firms.

# Bad luck

► The likelihood for the duration of job spells is independent of heterogeneity!

$$f(t) = \frac{\delta(\delta + \lambda)}{\lambda} \int_{\delta t}^{(\delta + \lambda)t} \frac{\exp(-x)}{x} dx.$$

► So we can identify $\lambda$ and $\delta$, and nothing about heterogeneity of firms and workers.

► (But the good thing is that we don't need to assume anything about it and we get $\delta$ and $\lambda$).

# Better luck

- Given bargaining on wages, outside options matter;
- and outside options generate option values, which increase with heterogeneity ( volatility!).
- "So" by looking at the distribution of wages we can infer heterogeneity.

# Application 2: moral hazard in insurance

Abbring-Chiappori-Pinquet, *JEEA* 2003.

- ▶ Insurees have exogenous types (risk) $v$ that are unobserved; we call this adverse selection;
- ▶ they also decide to adopt a risky behavior or not: moral hazard.
- ▶ Data typically gives us a series of claims for each individual.
- ▶ A state could be: "I have had exactly $p$ claims so far" and a spell is the time between two claims.

# Duration dependence

- Adverse selection induces positive duration dependence: the time between claims is positively correlated.
- On the other hand, with experience rating a claim (at fault) increases premia and makes risky behavior more costly—typically
- so moral hazard induces negative duration dependence.
- How can we test for the latter while controlling for the former?

# The Model

- The hazard function for claim $(p+1)$ at $t$, given state $p$, is (dropping $x$)

$$vh_0(t)A^{-p},$$

- with $A$ and $h_0$ unknown.
- $v$ models exogeneous unobserved risk,
- every time a claim occurs, the hazard for the next claim is divided by $A$: moral hazard.
- It is the MPH, with a twist: the $p$.

# Estimating Finite Mixtures

- In practice estimating finite mixture models can be tricky.
- A simple example is the mixture of normals (incomplete data likelihood)

$$f(x_1, \ldots, x_n | \theta) = \prod_{i=1}^{N} \sum_{k=1}^{K} \pi_k f(x_i | \mu_k, \sigma_k)$$

- We need to find both mixture weights $\pi_k = Pr(z_k)$ and the components $(\mu_k, \sigma_k)$ the weights define a valid probabiltiy measure $\sum_k \pi_k = 1$.
- Easy problem is label switching. Usually it helps to order the components by say decreasing $\pi_1 > \pi_2 > \ldots$ or $\mu_1 > \mu_2 > \ldots$
- The real problem is that which component you belong to is unobserved. We can add an extra indicator variable $z_{ik} \in \{0, 1\}$.
- We don't care about $z_{ik}$ per-se so they are nuisance parameters.

# Estimating Finite Mixtures

▶ We can write the complete data log-likelihood (as if we observed $z_{ik}$):

$$l(x_1, \ldots, x_n | \theta) = \sum_{i=1}^{N} \log \left( \sum_{k=1}^{K} I[z_i = k] \pi_k f(x_i \mu_k, \sigma_k) \right)$$

▶ We can instead maximized the expected log-likelihood where we take the expectation $E_{z|\theta}$

$$\alpha_{ik}(\theta) = Pr(z_{ik} = 1 | x_i, \theta) = \frac{f_k(x_i, z_k, \mu_k, \sigma_k) \pi_k}{\sum_{m=1}^{K} f_m(x_i, z_m, \mu_m, \sigma_m) \pi_m}$$

▶ Now we have a probability $\hat{\alpha}_{ik}$ that gives us the probability that $i$ came from component $k$. We also compute $\hat{\pi}_k = \frac{1}{N} \sum_{i=1}^{N} \alpha_{ik}$

# EM Algorithm

- Treat the $\hat{\alpha}_k(\theta^{(q)})$ as data and maximize to find $\mu_k, \sigma_k$ for each $k$

$$\hat{\theta}^{(q+1)} = \arg \max_\theta \sum_{i=1}^{N} \log \left( \sum_{k=1}^{K} \hat{\alpha}_k(\theta^{(q)}) f(x_i | z_{ik}, \theta) \right)$$

- We iterate between updating $\hat{\alpha}_k(\theta^{(q)})$ (E-step) and $\hat{\theta}^{(q+1)}$ (M-step)

- For the mixture of normals we can compute the M-step very easily:

$$
\begin{aligned}
\mu_k^{(q+1)} &= \frac{1}{N} \sum_{i=1}^{N} \hat{\alpha}_k(\theta^{(q)}) x_i \\
\sigma_k^{(q+1)} &= \frac{1}{N} \sum_{i=1}^{N} \hat{\alpha}_k(\theta^{(q)}) (x_i - \overline{x})^2
\end{aligned}
$$

# EM Algorithm

- EM algorithm has the advantage that it avoids complicated integrals in computing the expected log-likelihood over the missing data.

- For a large set of families it is proven to converge to the MLE

- That convergence is monotonic and linear. (Newton's method is quadratic)

- This means it can be slow, but sometimes $\nabla_\theta f(\cdot)$ is really complicated.

# My own example: Conlon Mortimer: AEJ 2013

- Probability of sales of $j$ depend on the set of available products $a_t$ some $x$'s (supressed) and some unknown parameters $\theta$ so that $p_j(a_t, \theta)$
- Imagine a multinomial logit with random coefficients
- At the beginning of the day/week we observe $a_t$.
- At some point during the week a product $k$ stocked out such that $a'_t = a_t \setminus k$.
- BUT we dont know which sales happened before or after the stockout.
- Now the probability of a sale is given by:
  $\lambda p_j(a_t, \theta) + (1 - \lambda) p_j(a'_t, \theta)$ where $\lambda$ is the fraction of consumers who arrive before the stockout.

# My own example: Conlon Mortimer: AEJ 2013

- Again we don't observe $(y_{jt}^0, y_{jt}^1)$ (sales before or after the stockout)
- However we do see the total sales $y_{jt} = y_{jt}^0 + y_{jt}^1$
- And we know something about product $k$ (the one that stocks out)
- We know that $y_{kt}^1 = 0$ by definition!
- We can compute the distribution of $f(\lambda|\theta, y_{kt})$ (when did the stockout occur?)

# My own example: Conlon Mortimer: AEJ 2013

What is $f(\lambda|\theta, y_{kt})$?

- How many consumers arrived before $y_{kt}$ consumers purchased good $k$?

- If sales are binomial then this is given by a <span style="color:red">negative binomial</span> distribution. This is an example of a <span style="color:red">hititng process</span> or <span style="color:red">discrete duration model</span>.

- Now I can compute

$$
\begin{aligned}
\hat{y}_{jt}^0 &= y_{jt} \int \frac{\lambda p_j(a_t, \theta)}{\lambda p_j(a_t, \theta) + (1 - \lambda) \cdot p_j(a_t', \theta)} f(\lambda|\theta, y_{kt}) d\lambda \\
\hat{y}_{jt}^1 &= y_{jt} - \hat{y}_{jt}^0
\end{aligned}
$$

- The M-step is just our usual MLE for logit, nested-logit, rc-logit treating $\hat{y}_{jt}$ as data.

# My own example: Conlon Mortimer: AEJ 2013

What was the point?

- We found that biases were large!
- Goods that stocked out a lot we understated demand for!
- Goods that were net recipients of substitution we overstated demand for!
- Basically you should pay attention to which other products are available.