

RAPID CALCULATION OF MEDICATION ADHERENCE USING PARALLEL COMPUTING WITH R AND PYTHON

Nick Davis, PhD

Assistant Professor of Research

Department of Medical Informatics

University of Oklahoma, Tulsa

School of Community Medicine



nicholas-davis@ouhsc.edu



[@argoneus](https://twitter.com/argoneus)

Oklahoma Supercomputing
Symposium

September 24, 2014

Medication Adherence

Definitions

- Adherence – fill prescriptions in a timely fashion
- Compliance – take medication as directed
- Persistence – deals with overall duration of drug therapy
- Subtle but distinct differences, but studies will sometimes use compliance and adherence interchangeably

Methodology

- Adherence is generally defined as the rate at which patients fill their prescriptions as indicated
- Drug claims data is considered to be an efficient and generally accurate means of assessing medication adherence
- One caveat is claims data is not a perfect representation of whether a patient is actually taking the medication
 - Patient may fill prescription and still fail to take medication

Methodology (cont.)

- Adherence generally expressed as a percentage
- Can be viewed roughly as how often/frequently patients take medication as directed
- Most studies have suggested a threshold of 80%
- Specific conditions may require a higher cutoff (e.g. 95%)

Methodology (cont.)

- Nonadherence is common for cardiovascular disease patients
- Psychiatric illness patients struggle with adherence, but have the greatest potential benefit – 58% among patients with psychoses
- *“Since adherence is enhanced when patients are involved in medical decisions about their care and in monitoring their care, the traditional model of the authoritarian provider should be replaced by the more useful dynamic of shared decision making by the health care provider and the patient.”*

Motivation

- Medication adherence called the “next frontier” in healthcare quality improvement
- Non-adherence is related to greater morbidity and mortality in chronic disease
- Non-adherence estimated to increase healthcare costs by over \$170 billion annually in the US alone
- Patient treatment and economic considerations suggest that non-adherence is a serious health concern that encourages research with a goal of impacting health outcomes and treatment costs
- *“Drugs don’t work in patients who don’t take them”* – former US Surgeon General C. Everett Koop



Measurement

- Centers for Medicare and Medicaid (CMS) have a number of recommended measures
- As suggested in a recent Pharmacy Quality Alliance (PQA) report, the Proportion of Days Covered (PDC) is the recommended medication adherence measure for electronic pharmacy claims
- PDC is defined as
$$\frac{\text{Number of days in period "covered" by medication}}{\text{Number of days in period}}$$
- PDC is a more conservative estimate of adherence than related measures in dealing with patients switching drugs

Data characteristics

- For this study, electronic pharmacy claims data from the Oklahoma Health Care Authority (OHCA) is used
- OU Tulsa receives a monthly pharmacy Medicaid claims feed from OHCA
- > 3.9 million claims records currently in data set
- Duration of claims is 1994 – present
- >97% of the data is from 2009 onward
- Over 134,000 unique Oklahoma Medicaid patients represented
- Over 20,000 unique medications represented

Data elements

Patient ID	Drug name
Date of Service (filled by patient)	Drug code
Date Prescribed	Drug class
DOB/Age	Drug quantity
Sex	Days supplied
Race	ProviderName (first, middle, last)
County	CoPay

Additional data sources

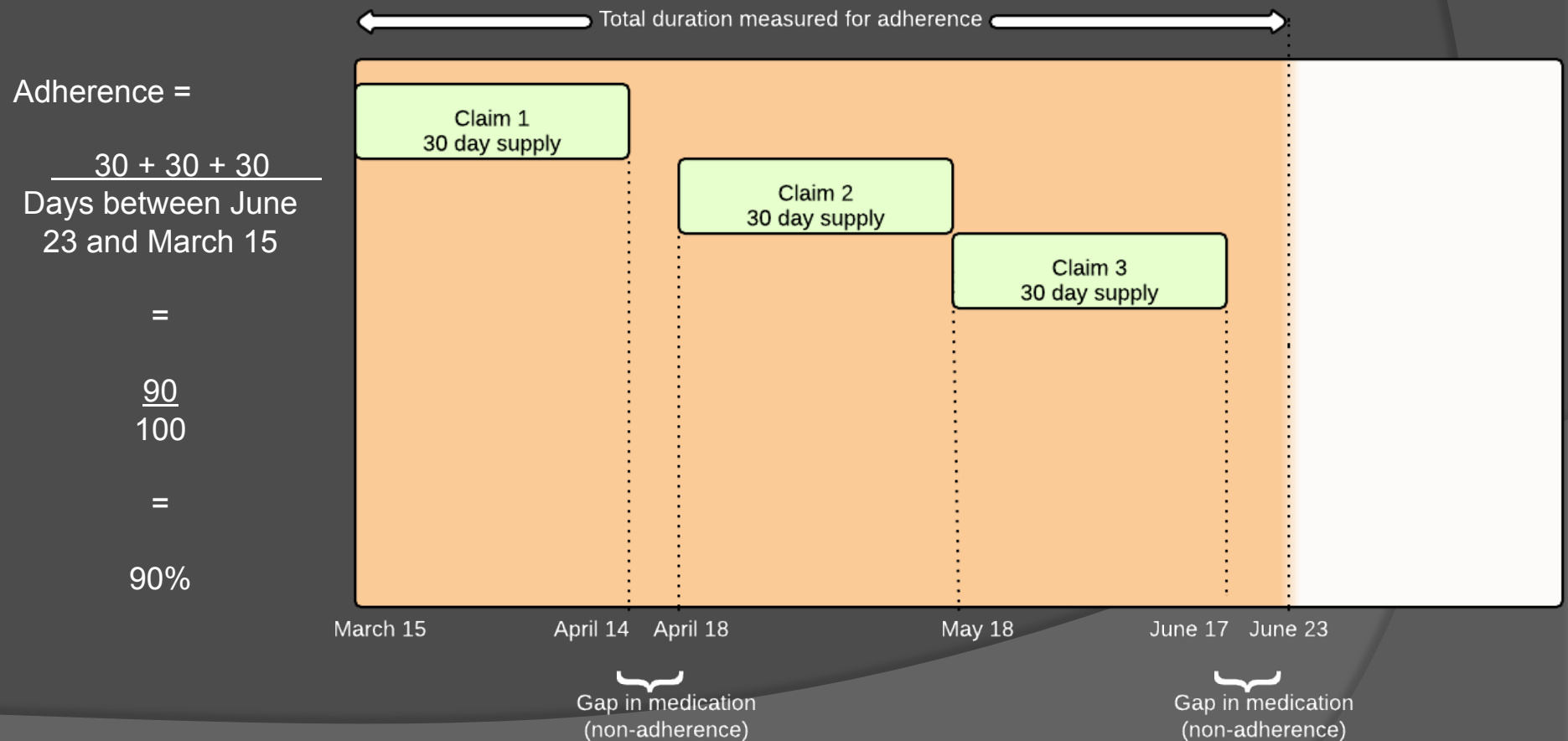
- NPI DB from CMS
- RxNorm <-> National Drug Code (NDC) crosswalk tables from UMLS
- Both have been loaded in DBs in the BI data warehouse
- Capture Provider details and RxNorm, VA Drug Class for most records

OHCA Data preparation

- Data aggregated by patient and medication combination
- The study period/duration is one year
- The total number of days covered by medication is determined by the number of days supplied for each prescription
- Key advantage:
Data is “embarrassingly parallel”

Example calculation

- Mitch has a prescription of a renin-angiotensin inhibitor for a cardiovascular condition. Each Rx is a 30 day supply (one pill per day), with refills provided on a monthly basis.

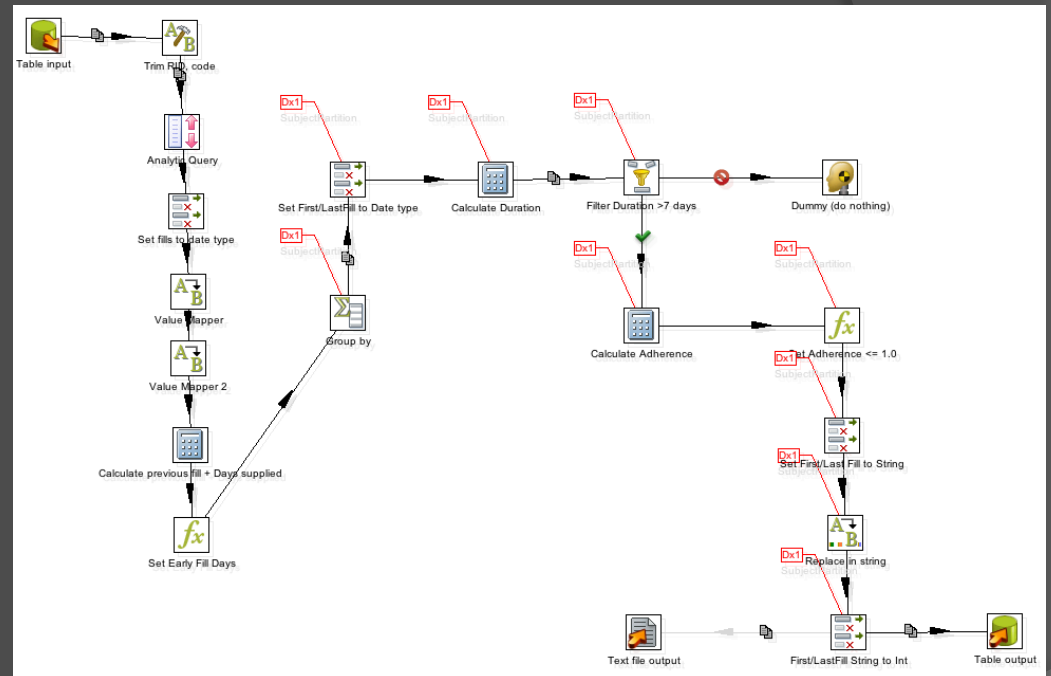


Testbed

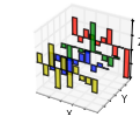
- ◎ Dell Precision T7610
- ◎ Intel Xeon E5-2697 2.7GHz, 12 cores
- ◎ 128 GB DDR3 RAM
- ◎ OS - Ubuntu Linux 14.04, kernel 3.13.0-24
- ◎ 100Mb/s Ethernet

Language choice

- OUSCM analytics environment leverages Pentaho
 - + Simple GUI-based programming
 - + Java based so relatively fast execution
 - Difficult to extend or modify, not a “real” programming language
 - Not widely used
- R and Python
 - + R is a premier language and platform for data analysis/data science applications
 - + Python has the Pandas library to provide R-like data structures (also Numpy/Scipy)
 - + Both are open source, large communities, lots of examples, libraries, documentation
 - Interpreted languages, slower than compiled programs



pandas
 $y_{it} = \beta x_{it} + \mu_i + \epsilon_{it}$



R Code

```
75 ▾ calc_adherence <- function(fills, year){
76   filldates <- as.Date(as.character(fills$FirstDateofService), format("%Y%m%d"))
77   days_supplied <- fills$DaysSupplied
78   first_fill <- min(filldates)
79   last_fill <- max(filldates)
80
81   duration <- as.Date(paste(year,"1231", sep=""), format("%Y%m%d")) - first_fill + 1
82   med_days <- vector(mode = 'integer', length=as.numeric(duration))
83 ▾   for (i in seq(duration)){
84 ▾     for (j in seq(length(filldates))){
85       if (filldates[j] <= first_fill + i - 1 && first_fill + i - 1 <= filldates[j] + days_supplied[j] - 1)
86         med_days[i] <- 1
87     }
88   }
89
90   early_fill_days <- calc_early_fill(filldates, days_supplied)
91   days_covered <- sum(med_days) + early_fill_days
92   adh <- days_covered / as.numeric(duration)
93   if (adh > 1) adh <- 1.0
94   ad <- data.frame(PatientID=fills$PatientID[1], DrugCode=fills$DrugCode[1], DrugName=fills$DrugName[1],
95                   DrugStrength=fills$DrugStrength[1], OHCALabel=fills$OHCALabel[1],
96                   OHCADrugClassName=fills$OHCADrugClassName[1], VADrugClassName=fills$VADrugClassName[1],
97                   FirstFill=first_fill,
98                   LastFill=last_fill, Duration=duration, DaysCovered=days_covered, MedicationDays=sum(med_days),
99                   Sex=fills$RecipientSexCode[1], Age=fills$Age[1], Race=fills$RecipientRaceCode[1],
100                  County=fills$RecipientCountyCode[1], CoPay=fills$CoPay[1], Zip=fills$Zip[1],
101                  DrugClass=fills$DrugTherapyClass[1], PrescribingProviderNPI=fills$PrescriberPhysicianProviderNPI[1],
102                  ProviderFirstName=fills$ProviderFirstName[1], ProviderLastName=fills$ProviderLastName[1],
103                  DrugQuantity=fills$DrugQuantity[1], LastDaysSupplied=tail(days_supplied, n=1),
104                  TotalDaysSupplied=sum(fills$DaysSupplied), EarlyDays=early_fill_days, Year=as.integer(year), Method="PDC",
105                  Adherence=adh)
106 }
```

R Code (cont.)

```
119 # split population by PatientIDs to partition the data for parallel processing
120 system.time(ad_pop_idx <- split(seq_len(nrow(ad_pop)), ad_pop$groupID))
121
122 system.time(par_adh <- mclapply(ad_pop_idx, function(i, ap, fun) fun(ap[i,], "2013"), ad_pop, calc_adherence, mc.cores = 1))
123
124 # merge list of data frames into a single data frame
125 system.time(adh_all <- rbindlist(par_adh))
126
127 completeFun <- function(data, desiredCols) {
128   completeVec <- complete.cases(data[,desiredCols])
129   data[completeVec, ]
130 }
131
132 # remove NAs for subjects with strange fill data (less than 7 days covered, leading to infinite adherence measures and
133 # other similar strangeness)
134 adh_filt <- completeFun(data.frame(adh_all), "Adherence")
135 # calculate mean adherence for all patients
136 mean(adh_filt$Adherence)
```


R Code (cont.)

```
119 # split population by PatientIDs to partition the data for parallel processing
120 system.time(ad_pop_idx <- split(seq_len(nrow(ad_pop)), ad_pop$groupID))
121
122 system.time(par_adh <- mclapply(ad_pop_idx, function(i, ap, fun) fun(ap[i,], "2013"), ad_pop, calc_adherence, mc.cores = 1))
123
124 # merge list of data frames into a single data frame
125 system.time(adh_all <- rbindlist(par_adh))
126
127 completeFun <- function(data, desiredCols) {
128   completeVec <- complete.cases(data[,desiredCols])
129   data[completeVec, ]
130 }
131
132 # remove NAs for subjects with strange fill data (less than 7 days covered, leading to infinite adherence measures and
133 # other similar strangeness)
134 adh_filt <- completeFun(data.frame(adh_all), "Adherence")
135 # calculate mean adherence for all patients
136 mean(adh_filt$Adherence)
```

Python Code

```
89 def calc_adherence(fills, year = None):
90     filldates = [ datetime.datetime.strptime(str(fills['FirstDateofService'][i]), '%Y%m%d').date()
91                 for i in fills['FirstDateofService'].index.values.tolist() ]
92     days_supplied = fills['DaysSupplied']
93     first_fill = min(filldates)
94     last_fill = max(filldates)
95     duration = ((datetime.date(int(year), 12, 31) if year else last_fill) - \
96                first_fill).days + 1
97     med_days = [0] * duration
98     for i in range(duration):
99         for j in range(len(filldates) if year else len(filldates) - 1):
100             if filldates[j] <= first_fill + datetime.timedelta(days=i) \
101                <= filldates[j] + \
102                datetime.timedelta(days=int(days_supplied.irow(j))) - datetime.timedelta(days=1):
103                 med_days[i] = 1
104     early_fill_days = calc_early_fill(filldates, days_supplied)
105     days_covered = sum(med_days) + early_fill_days
106     adh = days_covered / float(duration)
107     if adh > 1:
108         adh = 1.0
109     ad = pd.DataFrame({'PatientID' : fills['PatientID'][0],
110                      'DrugCode' : fills['DrugCode'][0],
111                      'RxNorm' : fills['RxNorm'][0],
112                      'RxNormLabel' : fills['RxNormLabel'][0],
113                      'DrugName' : fills['DrugName'][0],
114                      'DrugStrength' : fills['DrugStrength'][0],
115                      'OHCALabel' : fills['OHCALabel'][0],
116                      'OHCADrugClassName' : fills['OHCADrugClassName'][0],
117                      'VADrugClassName' : fills['VADrugClassName'][0],
118                      'FirstFill' : first_fill,
119                      'LastFill' : last_fill,
120                      'Duration' : duration,
121                      'DaysCovered' : days_covered,
122                      'MedicationDays' : sum(med_days),
123                      'Sex' : fills['RecipientSexCode'][0],
124                      'Age' : fills['Age'][0],
125                      'Race' : fills['RecipientRaceCode'][0],
126                      'County' : fills['RecipientCountyCode'][0],
127                      'CoPay' : fills['CoPay'][0],
128                      'Zip' : fills['Zip'][0],
129                      'DrugClass' : fills['DrugTherapyClass'][0],
130                      'PrescribingProviderNPI' : fills['PrescriberPhysicianProviderNPI'][0],
131                      'ProviderFirstName' : fills['ProviderFirstName'][0],
132                      'ProviderLastName' : fills['ProviderLastName'][0],
133                      'DrugQuantity' : fills['DrugQuantity'][0],
134                      'LastDaysSupplied' : days_supplied.irow(-1),
135                      'TotalDaysSupplied' : sum(fills['DaysSupplied']),
136                      'EarlyDays' : early_fill_days,
137                      'Year': int(year) if year != None else None,
138                      'Method': "PDC" if year != None else "smPDC",
139                      'Adherence' : adh }, index = [fills['groupID'][0]])
140     return ad
```

Python Code (cont.)

```
189 def run_analysis(year = '2013', pdc = True):
190     """Execute a series of steps required for calculating medication adherence"""
191     ad_pop = getRecords(True, year)
192     # set data type of groupID to int64
193     ad_pop['groupID'] = ad_pop['groupID'].astype(int)
194     # split population by groupIDs to partition the data for parallel processing
195     ad_pop_idx = ad_pop.groupby('groupID')
196     # Run the adherence calculation in parallel for each partition of data
197     start = time.clock()
198     par_adh = Parallel(n_jobs=-1)(delayed(calc_adherence)(
199         ad_pop.iloc[ad_pop_idx.groups[i]].reset_index(), (year if pdc else None)) for i in ad_pop_idx.groups)
200     end = time.clock()
201     print("calc: %f"%(end - start))
202     # merge list of data frames into a single data frame
203     start = time.clock()
204     adh_all = pd.concat(par_adh)
205     end = time.clock()
206     print("concat: %f"%(end - start))
207     # remove rows with NAs in Adherence column
208     adh = adh_all.dropna(axis=0, subset=['Adherence'])
209     # calculate mean adherence across all patients
210     print("mean adherence: %f"%adh['Adherence'].mean())
211     return adh
```

Python Code (cont.)

```
189 def run_analysis(year = '2013', pdc = True):
190     """Execute a series of steps required for calculating medication adherence"""
191     ad_pop = getRecords(True, year)
192     # set data type of groupID to int64
193     ad_pop['groupID'] = ad_pop['groupID'].astype(int)
194     # split population by groupIDs to partition the data for parallel processing
195     ad_pop_idx = ad_pop.groupby('groupID')
196     # Run the adherence calculation in parallel for each partition of data
197     start = time.clock()
198     par_adh = Parallel(n_jobs=-1)(delayed(calc_adherence)(
199         ad_pop.iloc[ad_pop_idx.groups[i]].reset_index(), (year if pdc else None)) for i in ad_pop_idx.groups)
200     end = time.clock()
201     print("calc: %f"%(end - start))
202     # merge list of data frames into a single data frame
203     start = time.clock()
204     adh_all = pd.concat(par_adh)
205     end = time.clock()
206     print("concat: %f"%(end - start))
207     # remove rows with NAs in Adherence column
208     adh = adh_all.dropna(axis=0, subset=['Adherence'])
209     # calculate mean adherence across all patients
210     print("mean adherence: %f"%adh['Adherence'].mean())
211     return adh
```

Runtimes

- 219K rows of data
- Each row represents a unique Patient/Medication/Fill date combination
- Reported times represent the average of 3 runs

Language	Number of cores	Overall Runtime	Calculation Runtime
R	1 (sequential)	1674.5 s	1652.7 s
Python	1 (sequential)	420.1 s	405.6 s
R	12	821.7 s	799.8 s
Python	12	132.9 s	115.6 s

Conclusions/Questions?

- R and Python (with Pandas) are excellent languages for data analysis
- Parallelizing code is often trivial, with some caveats
- Faster runtimes lead to richer exploration of the data