

教育・研究のためのデータ連携ワークショップ
医療・健康データ中心科学の実際
データクレンジングとマイニング

倉本 秋・片岡 浩巳
高知大学医学部

医療・健康データの発生元と統合DWH構築の課題



施設間誤差
測定方法間誤差
検体輸送時の誤差
検体変質の誤差



健康分野

医療分野

健診センター

複数の外部分析施設

病院、診療所

個人識別は数年後

特定機能健診

すでにオンライン収集
手順が確立され実施さ
れている分野
しかし、これらは連携さ
れていない、個別のプ
ロジェクト

院内感染
サーベランス

DPC関連データ

特定保健指導

個人識別ばらばら

診療データ

施設間誤差
測定方法間誤差

まだ大部分のシステム間の
データ交換は、紙ベースの運用

統合は、これから

共通の電子カルテ

家庭

グループ毎に分散された
データベース

多くのバイアスを含むデータ

プロトコルの標準化は完了し
ているがデータの正規化の
課題は多い



個人 (モバイル血圧計など)

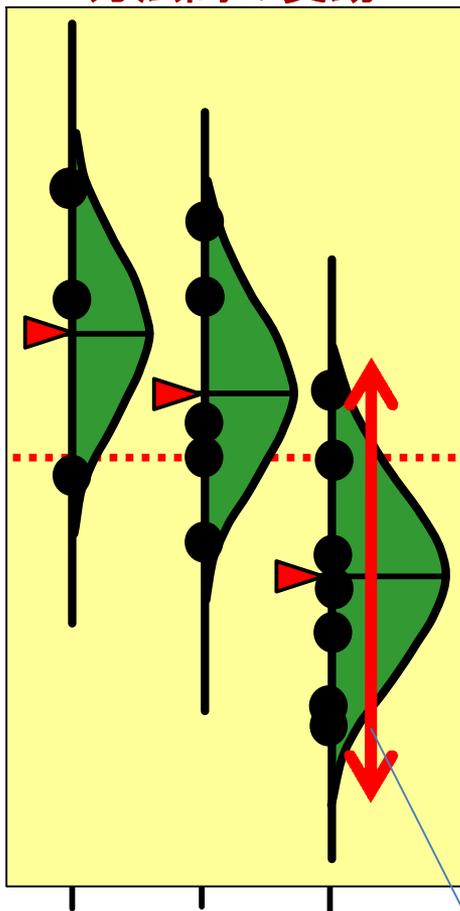
統合したデータウェアハウス

単にデータを集めただけの
DWHではマイニングは不能

方法間、施設間変動によるバイアスのクレンジング

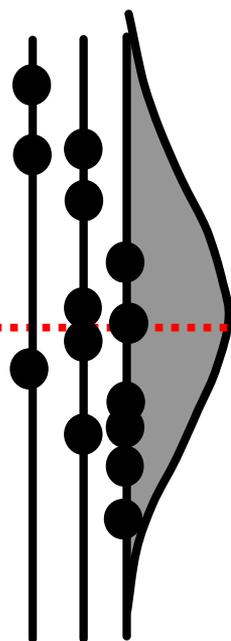
外部精度管理(同一の試料をすべての分析施設で測定した場合)

方法間の変動



施設間の変動
日差変動など

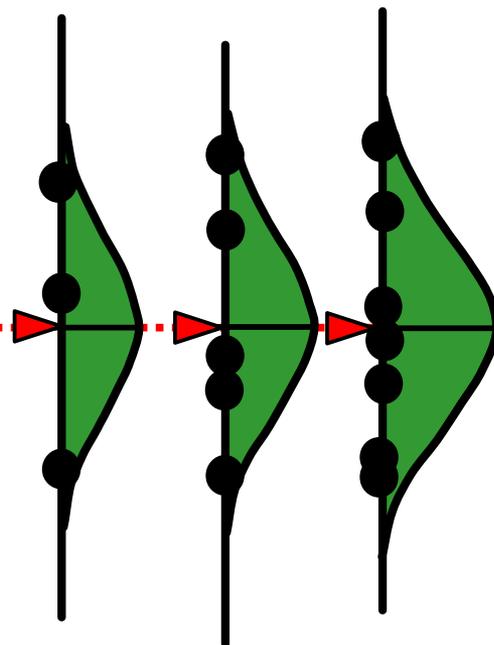
総変動



ただ集めただけのDWHでは
多くの誤差を含んだデータとなる

正規化

方法内変動

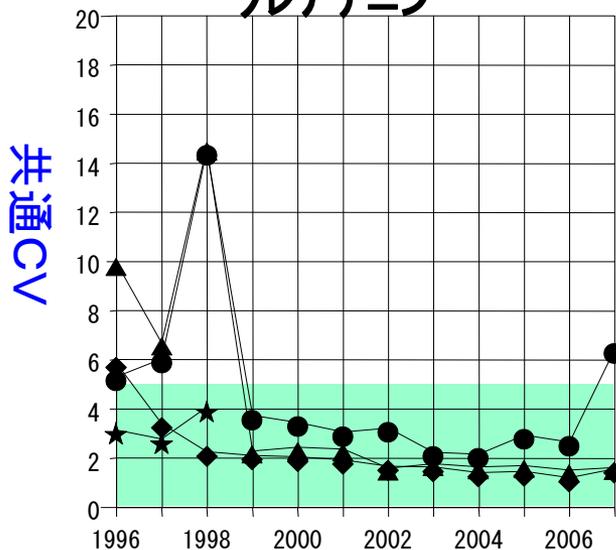


未来に向けては準備可能だが、これまで蓄積した過去の各病院の精度管理データの入手は困難

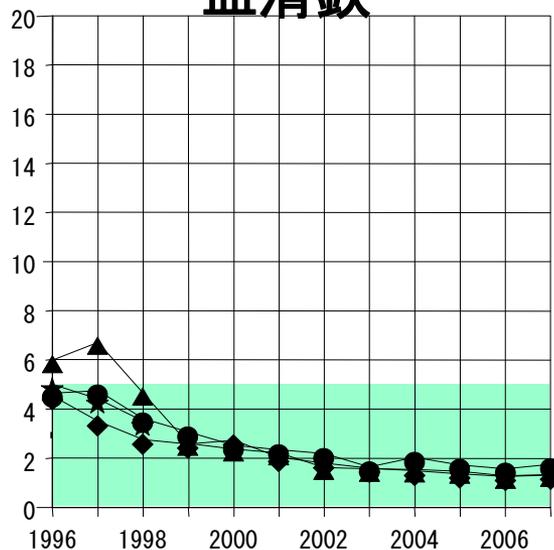
全国・外部精度管理の実態

過去と今

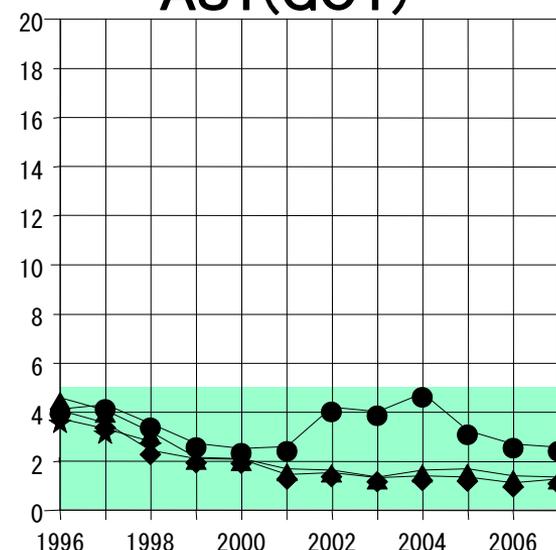
クレアチニン



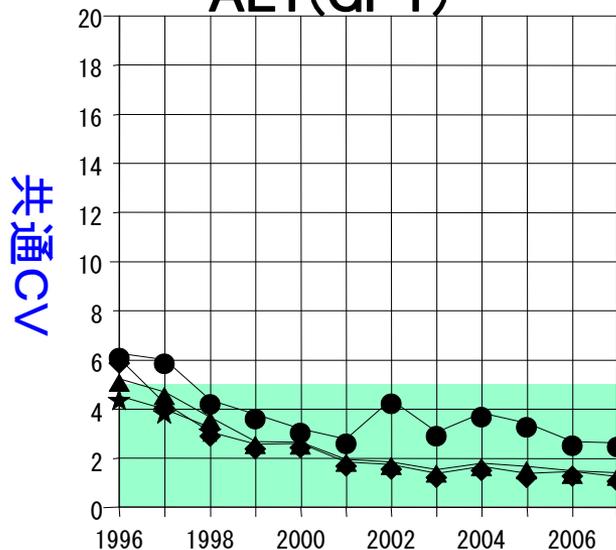
血清鉄



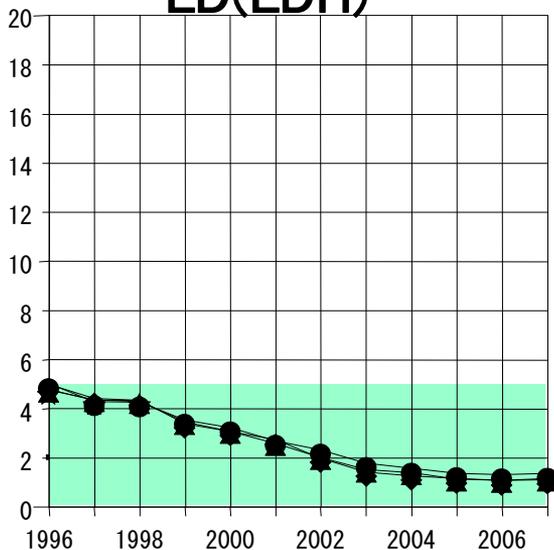
AST(GOT)



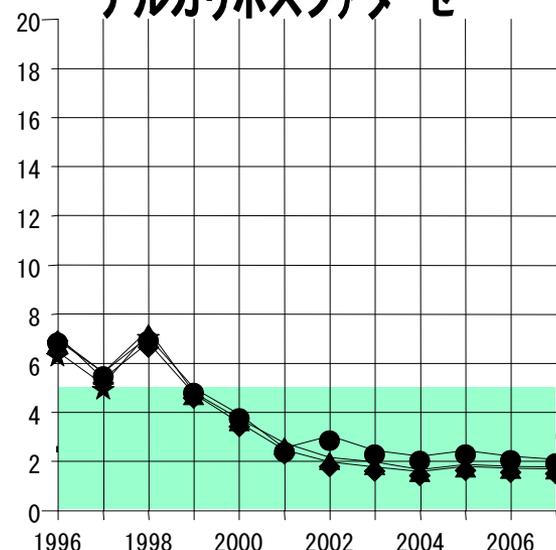
ALT(GPT)



LD(LDH)

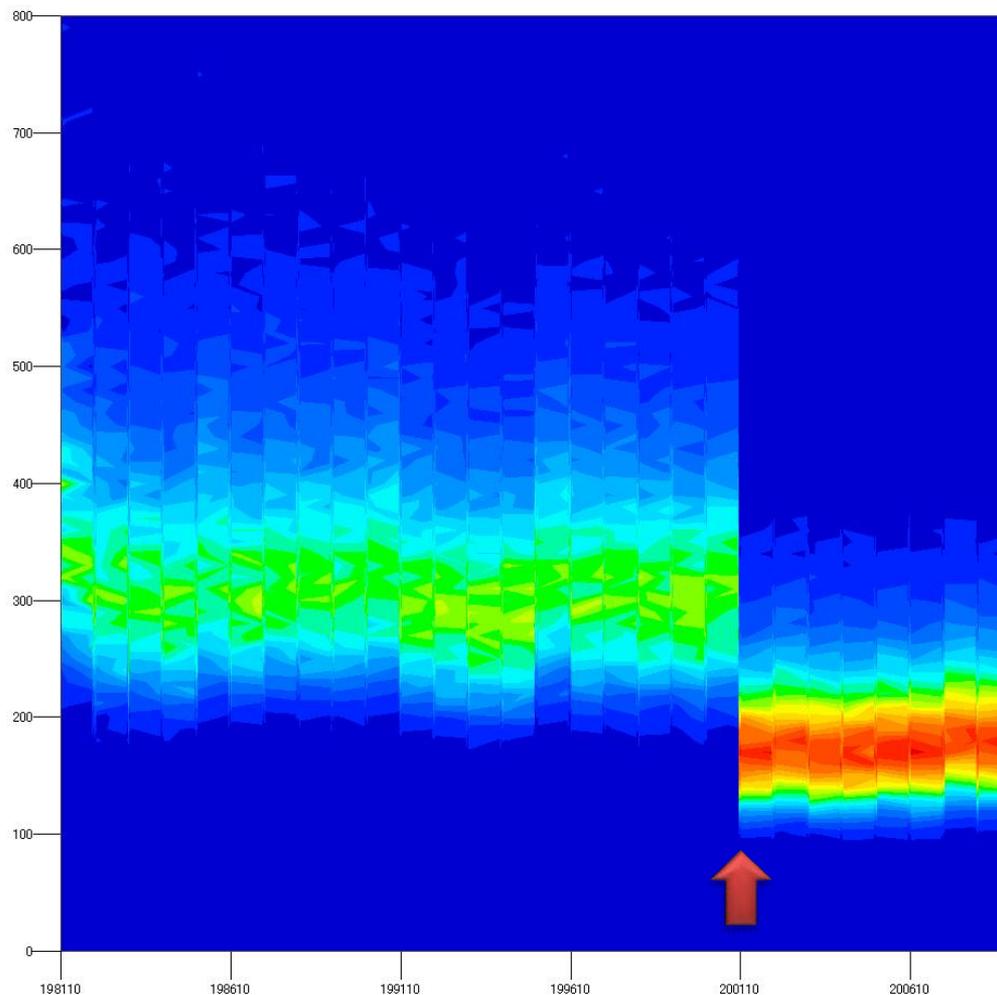
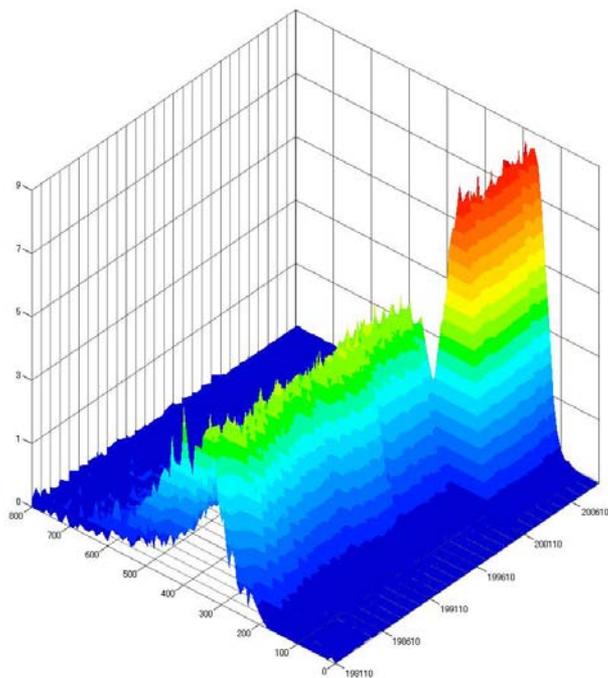


アルカリホスファターゼ



高知大学28年間の臨床検査値の分布

乳酸脱水素酵素(LD)の28年間にわたる集団分布(年月毎に度数分布を作成)



測定方法の変更

検査データの正規化技術の開発

目的: 蓄積された膨大なデータの分布から、変換のための係数を求める

長期の検査データ

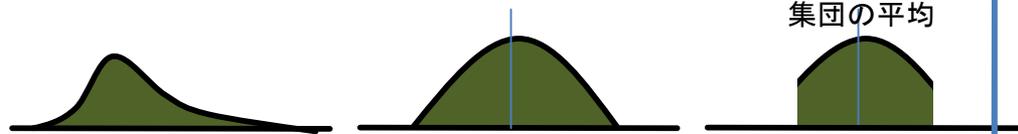


BoxCox変換

トランケーション

変換係数

旧測定法
1年毎の集団データ



検証用実験データ



旧測定法から新測定法への変換係数

回帰係数

ゴールドスタンダード

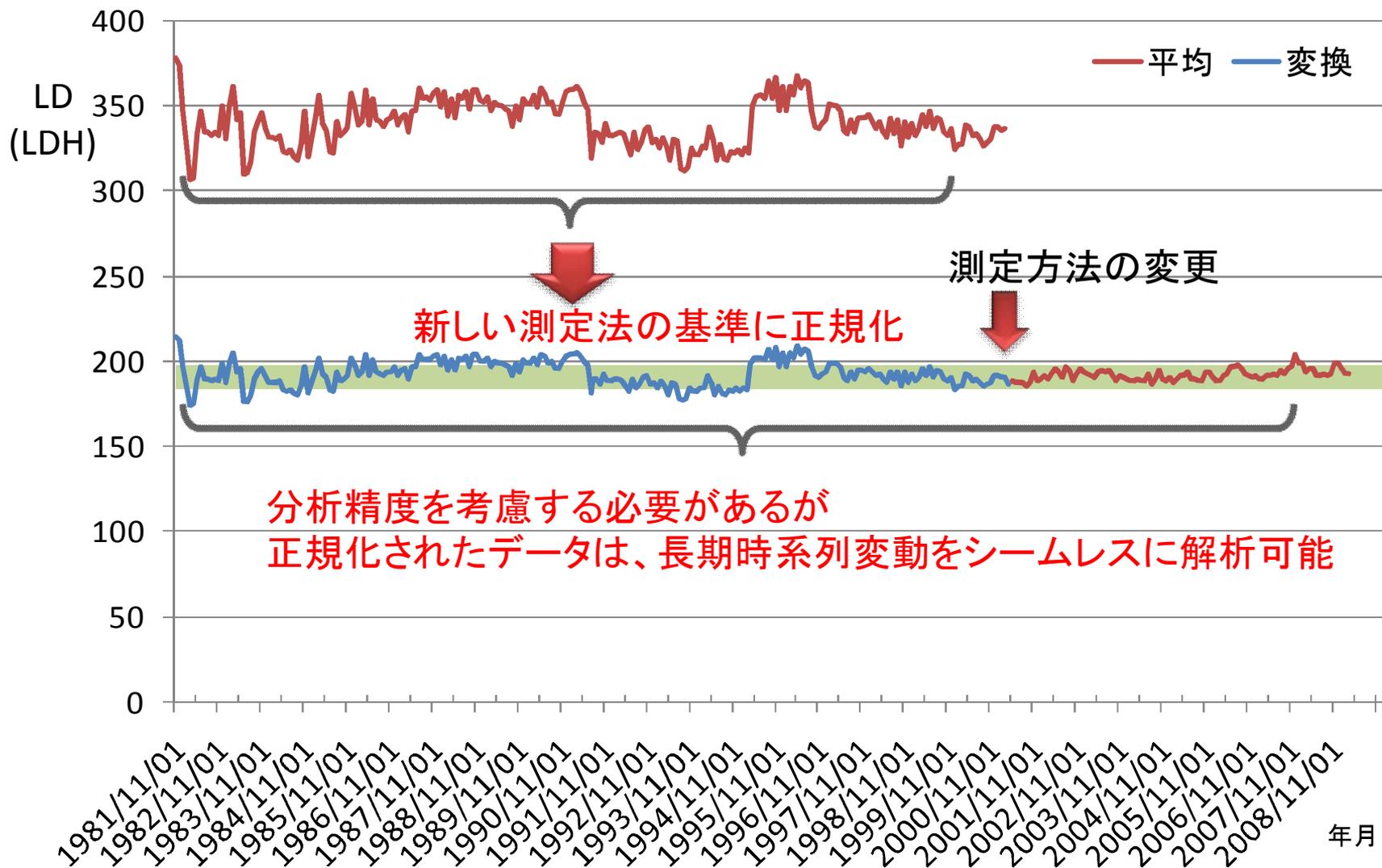
検証

一般的な検査室では、この実験データは無いし、過去の資料の記録が残っている施設は少ない

新旧測定法
同時測定データ

測定法の切り替え時に、新旧の測定法を使って数百件のサンプルを同時に分析する

正規化処理結果



まとめ

- 臨床検査データのクレンジング法の開発
 - 精度管理情報がなくても補正可能な手法
 - 蓄積されたデータの分布形を正規化し、その平均値を指標として変換
 - 実用レベルの成績が得られた
- 全国規模のDWH構築に向けての課題
 - 測定法、分析施設などの詳細情報を検査値と同時に収集、管理する必要がある