Wolfgang Lehner

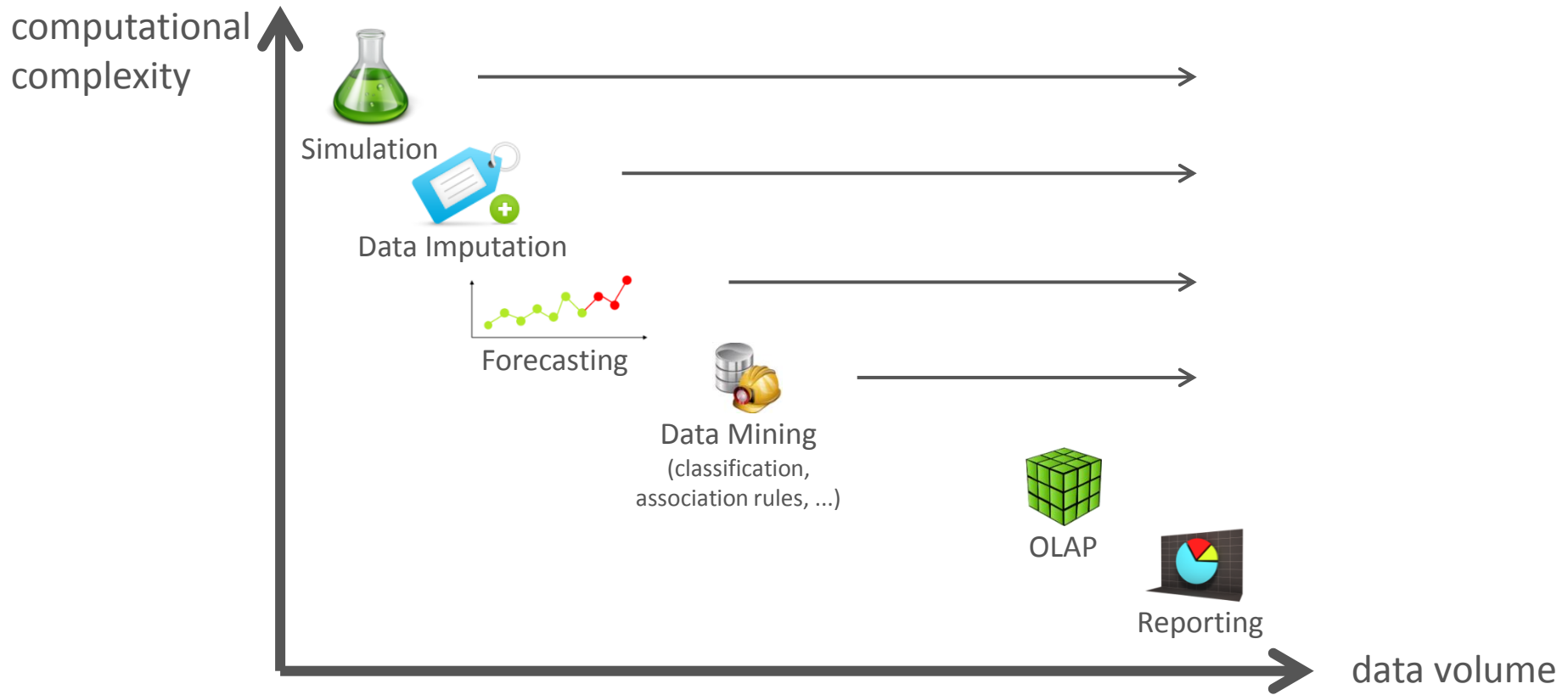# Forecasting and Data Imputation Strategies in Database Systems

11.07.2013

*data crunching* meets *number crunching*



computational complexity

Simulation

Data Imputation

Forecasting

Data Mining
(classification, association rules, ...)

OLAP

Reporting

data volume

BUSINESS

# The $1 Million Netflix Challenge

FRIDAY, OCTOBER 6, 2006 | BY KATE GREENE

*VP Jim Bennett discusses how recommendation systems suggest your next movie and the challenges of building a better one.*

✉ E-mail  ◁ Audio »  🖹 Print

Netflix' star rating system helps determine personalized movie recommendations. Now the company is looking to outside developers to improve those recommendations.

Earlier this week, Netflix, the online movie rental service, announced it will award $1 million to anyone who can come up with an algorithm that improves the accuracy of its movie recommendation service.

In doing so, the company is putting out a call to researchers who specialize in mac learning--the type of artificial intelligence used to build systems that recommend m books, and movies. The entrant who can increase the accuracy of the Netflix recommendation system, which is called Cinematch, by 10 percent by 2011 will w prize.

Recommendation systems such as those used by Netflix, Amazon, and other Wel retailers are based on the principle that if two people enjoy the same product, they' likely to have other favorites in common too.

But behind this simple premise is a complex algorithm that incorporates millions of ratings, tens of thousands of items, and ever-changing relationships between user preferences.
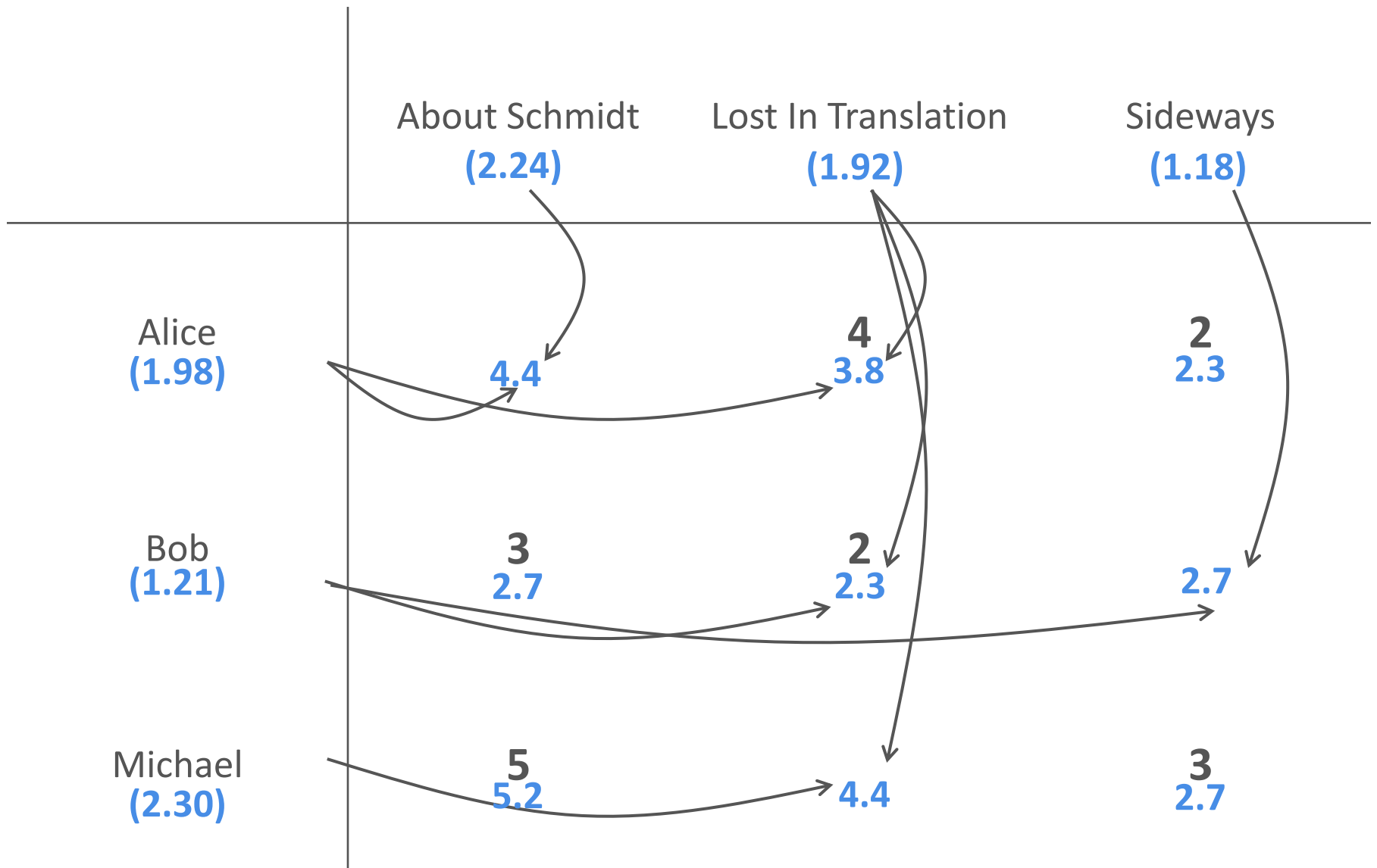
BellKor's Pragmatic Chaos

$$\hat{r}_{ui} = b_{ui} + |\mathrm{N}(u)|^{-\frac{1}{2}} \sum_{j \in \mathrm{N}(u)} e^{-\beta_u \cdot |t_{ui} - t_{uj}|} c_{ij} +$$

$$|\mathrm{R}(u)|^{-\frac{1}{2}} \sum_{j \in \mathrm{R}(u)} e^{-\beta_u \cdot |t_{ui} - t_{uj}|} \left( (r_{uj} - \tilde{b}_{uj}) w_{ij} \right) +$$
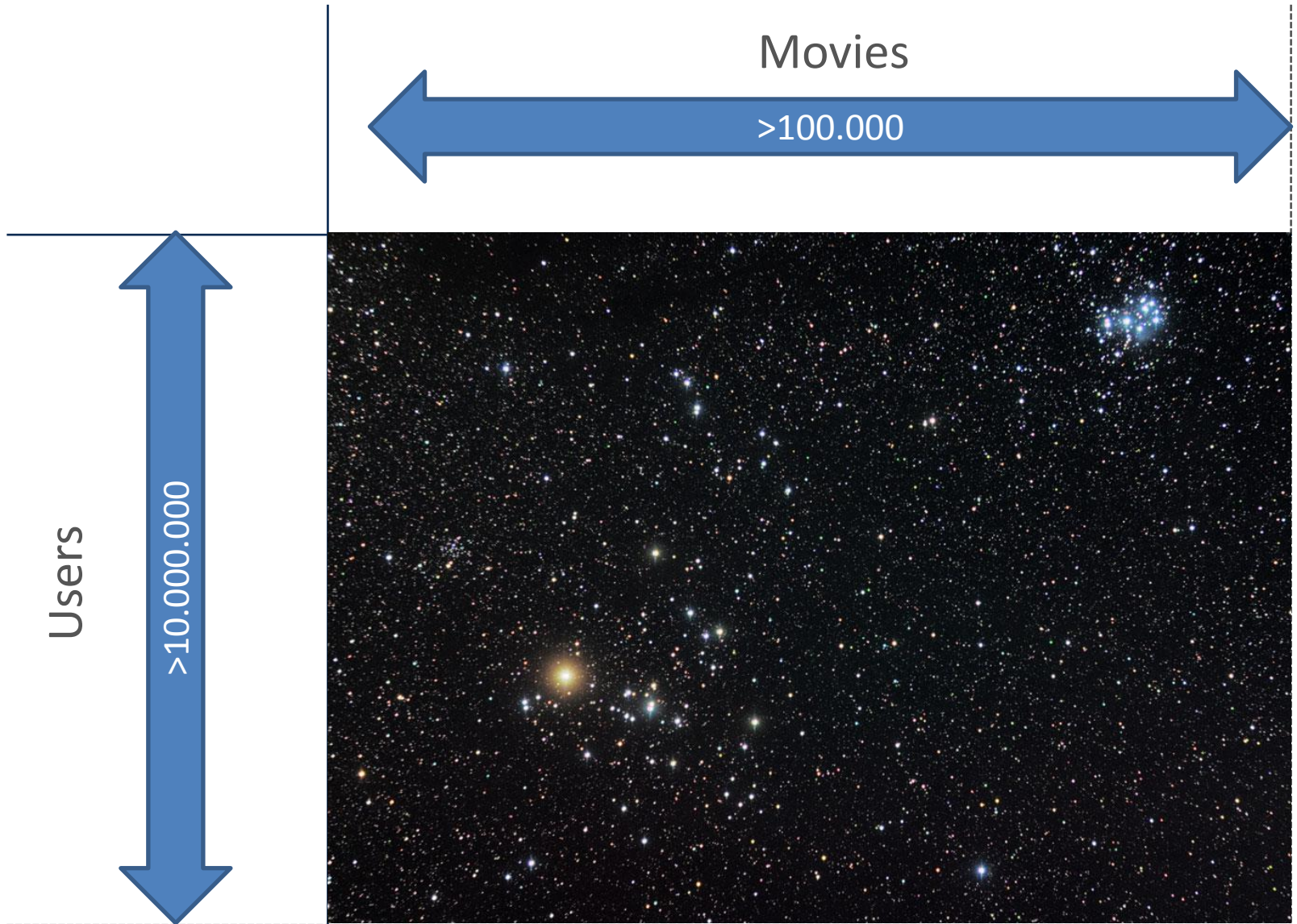
$$\sum_{j \in \mathrm{R}(u)} e^{-\gamma_u \cdot |t_{ui} - t_{uj}|} \left( (r_{uj} - \tilde{b}_{uj}) d_{ij} \right).$$

|  | About Schmidt (2.24) | Lost In Translation (1.92) | Sideways (1.18) |
|---|---|---|---|
| **Alice (1.98)** | **4** 4.4 | **4** 3.8 | **2** 2.3 |
| **Bob (1.21)** | **3** 2.7 | **2** 2.3 | 2.7 |
| **Michael (2.30)** | **5** 5.2 | 4.4 | **3** 2.7 |

Movies
>100.000

Users
>10.000.000

color code := user rating

different movies

different users

*Phase 1:* drop 75% of all pixels

**Phase 2:** Random permutation of rows and columns
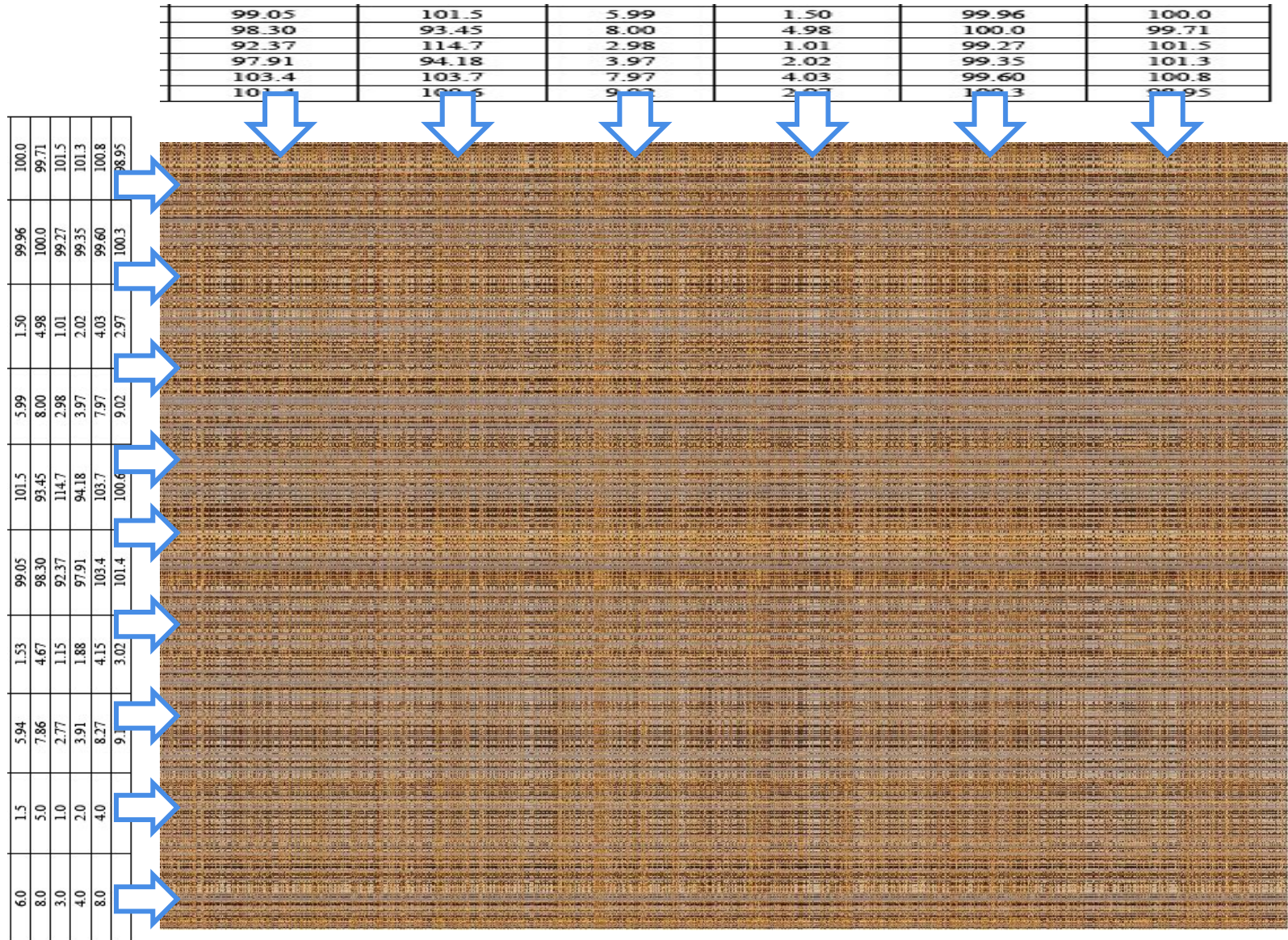
*Phase 3: Determine the latent factors*

Database Technology
Group

| | | | | | |
|---|---|---|---|---|---|
| 99.05 | 101.5 | 5.99 | 1.50 | 99.96 | 100.0 |
| 98.30 | 93.45 | 8.00 | 4.98 | 100.0 | 99.71 |
| 92.37 | 114.7 | 2.98 | 1.01 | 99.27 | 101.5 |
| 97.91 | 94.18 | 3.97 | 2.02 | 99.35 | 101.3 |
| 103.4 | 103.7 | 7.97 | 4.03 | 99.60 | 100.8 |
| 101.4 | 100.6 | 9.02 | 2.97 | 100.3 | 98.95 |

*Phase 4: Reconstruction*

| 99.05 | 101.5 | 5.99 | 1.50 | 99.96 | 100.0 |
| 98.30 | 93.45 | 8.00 | 4.98 | 100.0 | 99.71 |
| 92.37 | 114.7 | 2.98 | 1.01 | 99.27 | 101.5 |
| 97.91 | 94.18 | 3.97 | 2.02 | 99.35 | 101.3 |
| 103.4 | 103.7 | 7.97 | 4.03 | 99.60 | 100.8 |
| 101.4 | 100.6 | 9.02 | 2.97 | 100.3 | 98.95 |

*Phase 5: Final Result Generation*

25%

LFM

Distribution Grid

Central Power Plant

Network Access Point with Sensor

**Smart Grid**

Smart Homes

Smart Meters

Offices

Renewable Energy Sources

Data Management

Industrial Plant

Mobile Consumption Devices

Distributed Generation

Distributed Storage



➔ **Multi-Dimensional Time Series Data**

"It's tough to make predictions, especially about the future."
-- Mark Twain

*Given*

- Time series with numerical values as training data

*Goal*

- Predict future values for arbitrary future point in times (forecast horizon)
- Include trend and seasonality

*Applications*

- Planning of sales and budget
- Price development
- Inventory, manufacturing
- Climate, weather, environment
- Economic indicators
- Stocks

**sales (memory cards)**

*Diverse Time Series Data*

- Find forecast model that minimizes the forecast error

*Example: Exponential Smoothing Framework*

- 30 variants (additional additive or multiplicative errors)

| Trend/Season | No Season | Add. Season | Mult. Season |
|---|---|---|---|
| No Trend | | | |
| Add. Trend | | | |
| Add. Damp. Trend | | | |
| Mult. Trend | | | |
| Mult. Damp. Trend | | | |

*Mathematical Foundations*

*In-DBMS Time Series Forecasting*

- Flash-Forward Query Project

Flash Forward
Query

*Query Processing and Optimization*

*Model Configuration Advisor*

- Forecasting the Data Cube
- Selection of Model Configurations

*Notification-based Forecast Queries*

*Beyond Forecast Models*

- Model Refinement

# *Mathematical Foundations*

## Forecast Model

- Statistical time series description (model)

## (Recursive) Forecasting Process

- Model Identification
- Model Estimation
- Forecasting and Model Update
- Model Evaluation
- Model Adaptation

**Model Maintenance**

New Time Series Values $U_i$

Forecasting Values $F_{i+h}$

Model Parameters
$\phi_1=0.55$, $\phi_2=0.45$

**Forecasting**

**Model Evaluation**

**Time Series Data**

Model Type
**AR(2)**

**Model Estimation**

**Model Identification**

**Model Adaptation**

**Model Creation** | **Model Usage**

# > Model Identification

Database Technology Group

*Forecast Model Types / Classes*

**Domain-Specific Extensions** ← **Base Forecast Models**

White-Box

**(Auto)Regression**  **Exponential Smoothing**  **Machine Learning**

Gray-Box  Black-Box

**HWT**
(Single-Equation)

**EGRV**
(Multi-Equation)

**AR**
**MA**
**ARMA**
**ARIMA**
**SARIMA**
**ARMAX**

**MLR**
(Multiple Linear Regression)

**SESM**
(Single Exponential Smoothing)

**DESM**
(Double Exponential Smoothing)

**TESM / HoltWinters**
(Triple Exponential Smoothing)

**BN**
(Bayesian Networks)

**SVM**
(Support Vector Machines)

**SVR**
(Support Vector Regression)

**ANN**
(Artificial Neural Networks)

31211ok, enough. Let me finalize.

## Problem

- Instantiate a forecast model w.r.t. meta model and training data set

$e_{MSE}=827,354.4$

## Example

- Forecast Model Type **AR(2)**:

$$\hat{y}_t = \phi_1 \cdot y_{t-1} + \phi_2 \cdot y_{t-2}$$

- Error Metric: **MSE**

$$\frac{1}{n}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2$$

- Parameter Estimator **L-BFGS-B**



$e_{MSE}=$ 211,204.7

7 iterations ( 35 fn eval )

0.23

0.56

*Forecasting*

- Use the estimated forecast model   $\hat{y}_t = 0.56 \cdot y_{t-1} + 0.23 \cdot y_{t-2}$
- Create h forecast values (forecast horizon)
- Update model state for new measurements (e.g., exponential smoothing)

*Example Forecast (EGRV)*

- SMAPE$_{vshort}$ **=0.0021**
- SMAPE$_{long}$ **=0.0755**

## Model Evaluation

- Goal: Trigger model adaptation only if necessary

- Fixed Interval Techniques
  (# updates, time interval)
- Continuous Evaluation
  Techniques
  (threshold, on-demand)

h=2

## Model Adaptation

- Goal: Adapt the forecast model to the changed time series (if necessary)

- Model Re-Identification
- Model Re-Estimation (old model as start point)

# *In-DBMS Time Series Forecasting*

Flash Forward Query

## Magnetic

- **„Attract data and practitioners"**
- Use all available data sources independent of their quality

## Agile

- **„Rapid iteration: ingest, analyze, productionalize"**
- Continuous and rapid evolution of physical and logical structures
- ELT (Extraction, Loading, Transformation)

## Deep

- **„Sophisticated analytics in Big Data"**
- Extended algorithmic runtime environment
- Ad-hoc advanced analytics and statistics

1. **mad skills**                                    92 up,
   To be able to do/perform amazing/unexpected things
   *I gots me mad skills, yo.*
   *To be said after performing an extraordinairy feat.*

*Integration of advanced analytics into scalable database management systems*

- Traditionally forecasting is performed manually in external statistical systems
- Support of transparent and automatic in-DBMS forecasting

**Traditional**

**In-DBMS Forecasting**

**Analysis (e.g., R, SPSS)**

**DBMS**

**Predictive DBMS**

## Model-based Prediction

- We employ (classical) time series models (e.g. exponential smoothing, ARIMA)
- Used to
  - Explain time series from it's history
  - And—possibly—from exogenous inputs

## Related Work

- Customized functions with proprietary languages
  - SQL Server 2012: ARIMA, autoregressive trees
  - Oracle: exponential smoothing, non-linear regression
- Bi-directional communication
  - Reuse existing statistical tools (e.g. R)
  - SAP HANA, Oracle, IBM Netezza
- Model-based views

## Key Elements

- Declarative querying
- Automatic model creation
- Automatic model maintenance
- Forecast model advisor

t=6 ⟶ t=7

DWH

Model ⟶ Forecast

| 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 |

SELECT … FROM sales

SELECT … FROM sales  **AS OF**  …

DWH

Model ⟶ Forecast

| 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 |

SELECT … FROM sales

SELECT … FROM sales  **AS OF**  …

Evaluate Error

Adjust the Model

Forecast Queries          Inserts

**Query Interface**

**Transparent Model Usage**
- Model matching
- Reuse of models
- Optimization of adhoc-queries

Model Index

**Model Maintenance**
- Model matching
- Model evaluation
- Model maintenance

Model Pool

Model

Model

Model

Model

**Model Creation**
- Manually by experts
- Automatically
- Ensemble-based modeling

Time Series          Time Series          Time Series

**Model Advisor**
- Physical design recommendation

Base Tables

# *Query Processing and Optimization*

*SQL Forecast Query*

```
SELECT  Date, Quantity
  FROM  Sales, Article
  WHERE S_ID = A_ID AND
        Name = 'Article A'
  AS OF 2013-05-02
```

| Date | Qty |
|------|-----|
| 2013-04-28 | 6 |
| 2013-04-29 | 5 |
| 2013-04-30 | 7 |
| **2013-05-01** | **6** |
| **2013-05-02** | **6.5** |

*Query Plan*

- Forecast operator Ψ

$\Psi_{k=2}$

|

$\pi_{Date, QTY}$

|

$\bowtie_{S\_ID=A\_ID}$

Sales

$\sigma_{Name= \,'Article A'}$

| ... | S_ID | Date | Qty |
|-----|------|------|-----|
| ... | 1 | 2013-04-28 | 6 |
| ... | 1 | 2013-04-29 | 5 |
| ... | 2 | 2013-04-30 | 1 |
| ... | 1 | 2013-04-30 | 7 |

Article

| A_ID | Name | ... |
|------|------|-----|
| 1 | Article A | ... |
| 2 | Acticle B | ... |

*Selection invariance*

- Condition filtering out whole clusters

*Projection invariance*

- If neither time, measure or cluster attribute

*Union invariance*

- On different relations

$\Psi_{k=2}$

$\pi_{date, quantity}$

$\sigma_{product= 'HTC'}$

sales

$\Longleftrightarrow$

$\sigma_{product= 'HTC'}$

$\Psi_{k=2}$

$\pi_{product, date, quantity}$

sales

➔ Influence

| date | product | quantity |
|------|---------|----------|
| Jun 2013 | HTC | 36.000 |
| Jul 2013 | HTC | 38.000 |
| Jun 2013 | Nokia | 5.000 |
| Jul 2013 | Nokia | 3.000 |

$\cup$

$\Psi_{k=2}$        $\Psi_{k=2}$

$\sigma_{product= 'HTC'}$    $\sigma_{product= 'NOKIA'}$

sales1          sales2

$\Longleftrightarrow$

$\cup$

$\sigma_{product= 'HTC'}$    $\sigma_{product= 'NOKIA'}$

sales1          sales2

*Restructuring might also influence accuracy*

- E.g. join, aggregation

$\pi_{sales1.date,\ sales1.quantity\ -\ sales2.quantity}$

$\bowtie$ sales1.date=sales2.date

$\Psi_{k=2}$        $\Psi_{k=2}$

sales1      sales2

$\Psi_{k=2}$

$\pi_{sales1.date,\ sales1.quantity\ -\ sales2.quantity}$

$\bowtie$ sales1.date=sales2.date

sales1      sales2

Expensive model creation                Expensive join

➔ Cost model required: accuracy and runtime

*How to compute aggregation forecast queries?*

Base time series

| | Runtime | Accuracy |
|---|:---:|:---:|
| | ⬆ | ⬇ |
| | ⬇ | ⬆ |
| | ➡ | ⬆ |

|  | t-1 | t | t+1 |  |
|---|---|---|---|---|
| base time series | ... 8 | 9 | → 10 | Model |
|  | ... 48 | 49 |  |  |
|  | ... 8 | 9 | → 10 | Model |
| aggregate | ... 64 | 67 | → $\hat{y}?$ | $= \alpha \sum Forecast$ |

*Uniform Estimation*
$$\alpha = \frac{\# \, base \, series}{\# \, base \, series \, in \, sample} \qquad \frac{3}{2} \cdot 20 = 30$$

*Estimation with Historical Ratios*

$$ratio = \frac{base \, series}{aggregate \, series} \quad \Rightarrow \quad \alpha = \frac{1}{\sum ratios} \qquad \frac{1}{\frac{9}{67}+\frac{9}{67}} = 3.7 \cdot 20 = 74$$

*Calculation of Historical Ratios*

- Different approaches possible
    - Simple averages
    - Lagged proportions
    - …
- Seasonality of data is important

|  | t-1 | t | t+1 |
|------|-----|-----|-----|
| base … | 8 | 9 | 10 |
| agg … | 64 | 67 | |

*Combined strategy: mixture of past ratios and ratio one season ago*

time

Jul 11          Feb 12          Jul 12

*agg*



*base*

*Aggregation - Sales*

# *Model Configuration Advisor*

city → region

C1
C2
R1

C3
C4
R2

time

P1   P2   P3   P4

product

```
SELECT    time, measure
FROM      facts
WHERE     product = P4
AND       city = C4
AS OF     now() + 1 day
```

measure

time

city → region

C1
C2
R1

C3
C4
R2

```
SELECT    time, SUM(measure)
FROM      facts
WHERE     product = P4
AND       region = R2
GROUP BY  time
AS OF     now() + 1 day
```

time

P1    P2    P3    P4

product

M  M  M

*Traditional Database Systems*

*Statistical Database Systems*

*Conceptually, we organize the aggregation possibilities as a <u>directed time series hyper graph</u>*



**Aggregated Time Series**

**Base Time Series**

*Each node (or time series) may be associated with a forecast model*



*Time series methods*

- Exponential Smoothing
- ARIMA

*A query describes one or several nodes in the hyper graph*

*Forecast values of a node can be computed by any nodes in the graph*

*Derivation weight k*

- Based on history of source and target time series

**k**

**M**

Direct

**k=1**  **M**

Aggregation

**∑k=1**

**M**  **M**  **M**

Disaggregation

**k<1**  **M**

*Model configuration*

Assignment of models

**+**

Assignment of derivation schemes



*Configuration evaluation*

Costs        Forecast accuracy

*Heuristically indicate the expected benefit of a model at a node (without building the model)*

- Focus on time series relationships
- Measure to specify the <u>derivation error</u> between two nodes
- Low indicator value →low error (*good derivation*)
- High indicator value → high error (*poor derivation*)

| 10, 20, 32, 40 | | 10, 20, 32, 40 |
|---|---|---|
| | ◯ | |
| 0, 0, 0.06, 0 | **Historical Error** | **Weight Variance** 0.1, 0.1, 0.09, 0.1 |
| 1, 2, 3, 4 | 🔴 | 1, 2, 3, 4 |
| | **Indicator** | |

Database **Technology**
Group

*Local indicator arrays*

→ *Derivation errors of one node*

*Global indicator array*

→ *Minimum over all local arrays*



**MIN**

*Select start configuration*

*Select start configuration*

*Candidate selection*

- Preselection

*Select start configuration*

*Candidate selection*

- Preselection

M

| 0 | 5 | 2 | 3 |
|---|---|---|---|

Add    Delete

*Select start configuration*

*Candidate selection*

- Preselection
- Ranking

**M**

| 0 | 0 | 1 | 3 |
|---|---|---|---|

| 0 | 5 | 2 | 3 |
|---|---|---|---|

1    2

Add          Delete

*Select start configuration*

*Candidate selection*

- Preselection
- Ranking

*Evaluation*
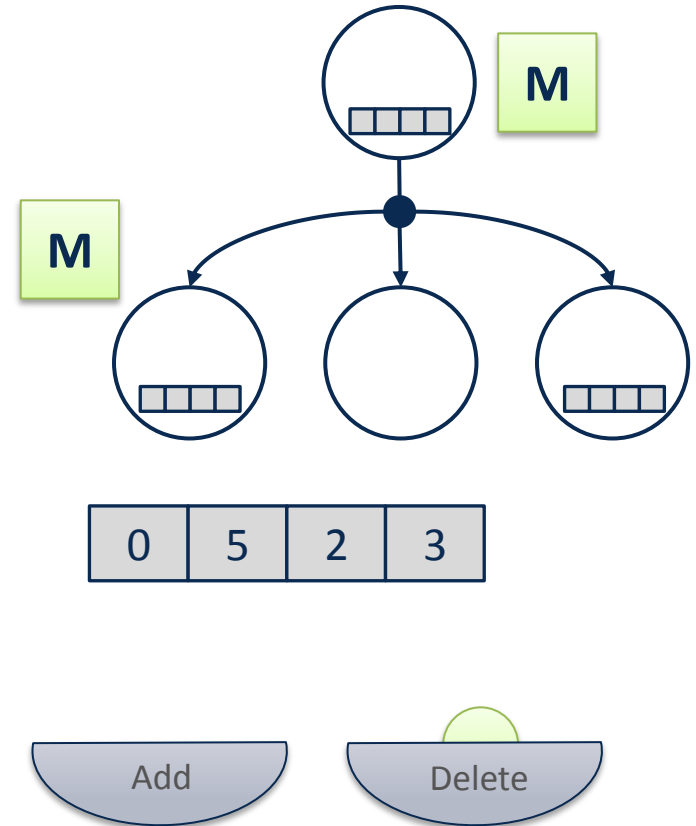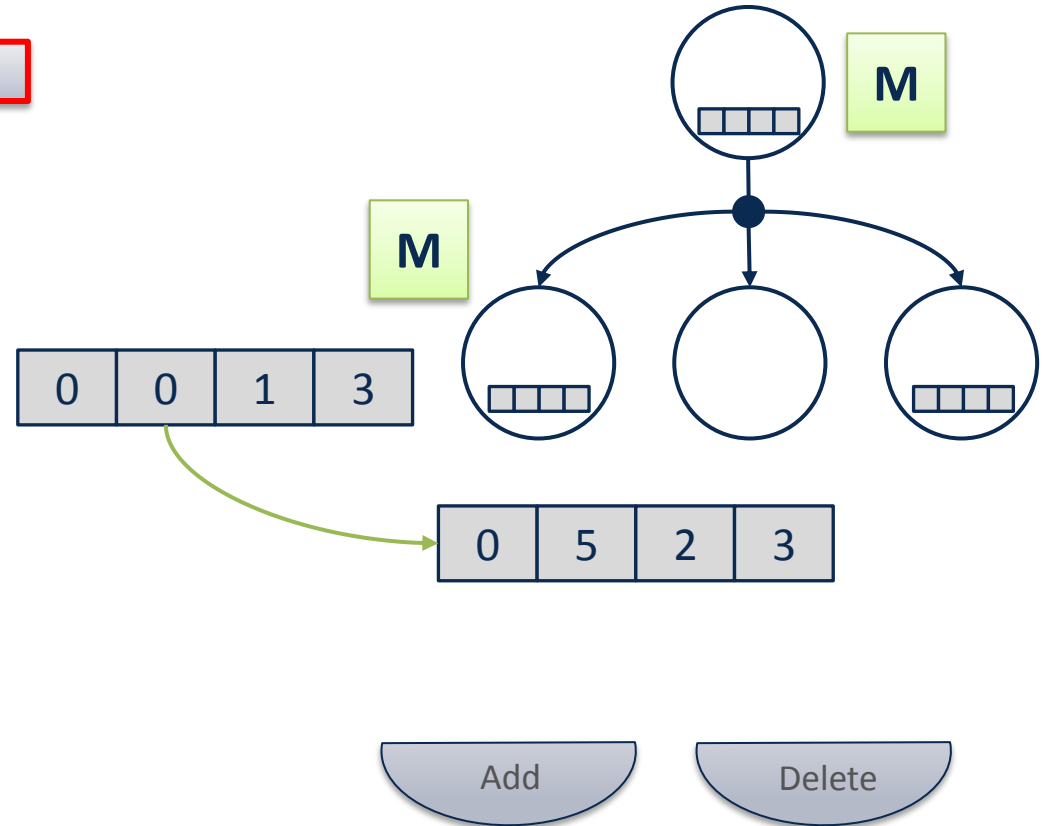
- Model creation
- Acceptance

*Select start configuration*
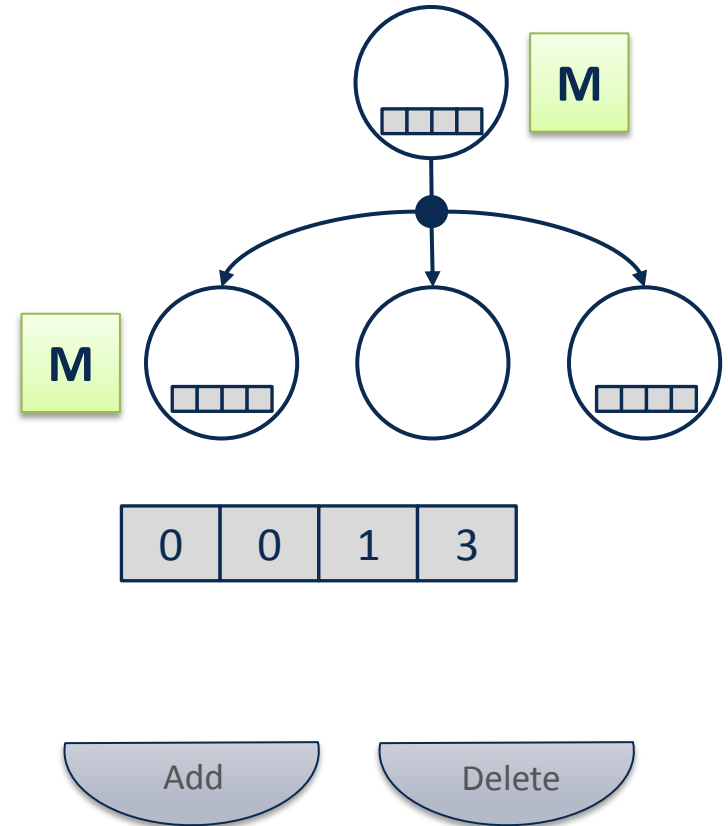
*Candidate selection*

- Preselection
- Ranking

*Evaluation*

- Model creation
- Acceptance

M

M

| 0 | 5 | 2 | 3 |
|---|---|---|---|

Add    Delete

*Select start configuration*

*Candidate selection*

- Preselection
- Ranking

*Evaluation*

- Model creation
- Acceptance

M

M

| 0 | 0 | 1 | 3 |
|---|---|---|---|

| 0 | 5 | 2 | 3 |
|---|---|---|---|

Add

Delete

Database **Technology**
Group



*Select start configuration*

*Candidate selection*

- Preselection
- Ranking

*Evaluation*

- Model creation
- Acceptance

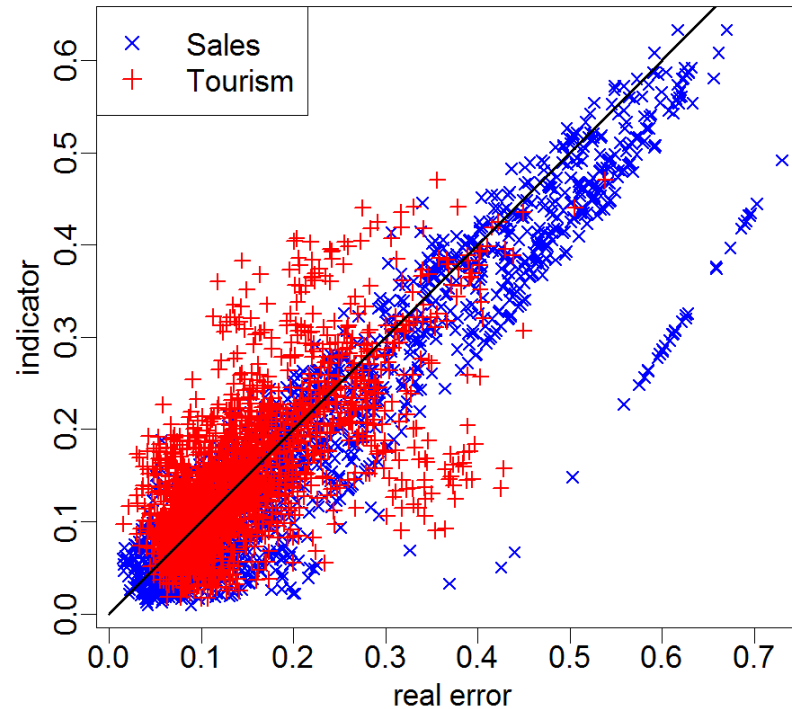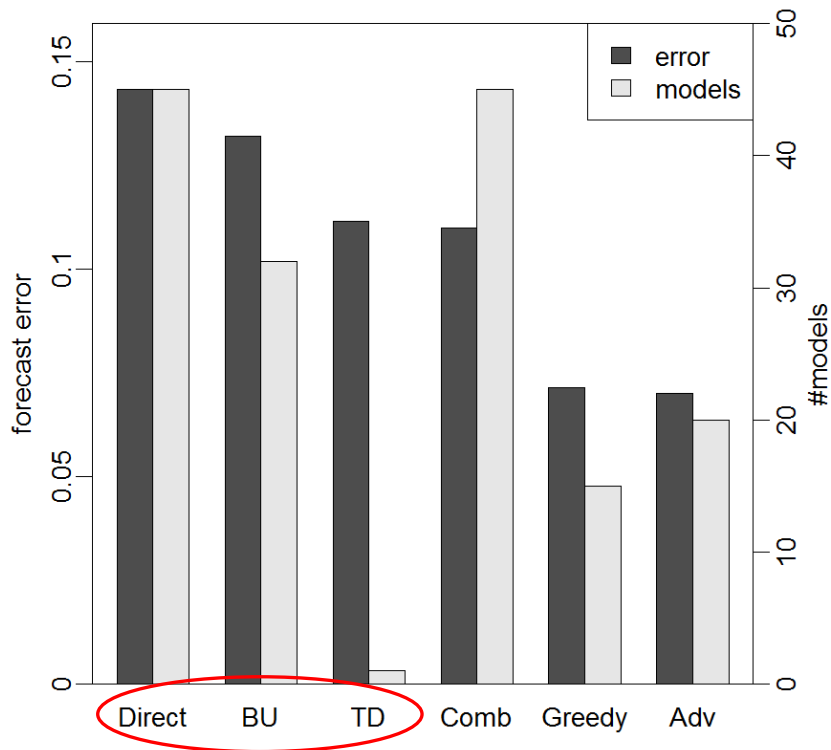| 0 | 0 | 1 | 3 |
|---|---|---|---|

Add      Delete

Model costs: 2      Forecast error: 10 %
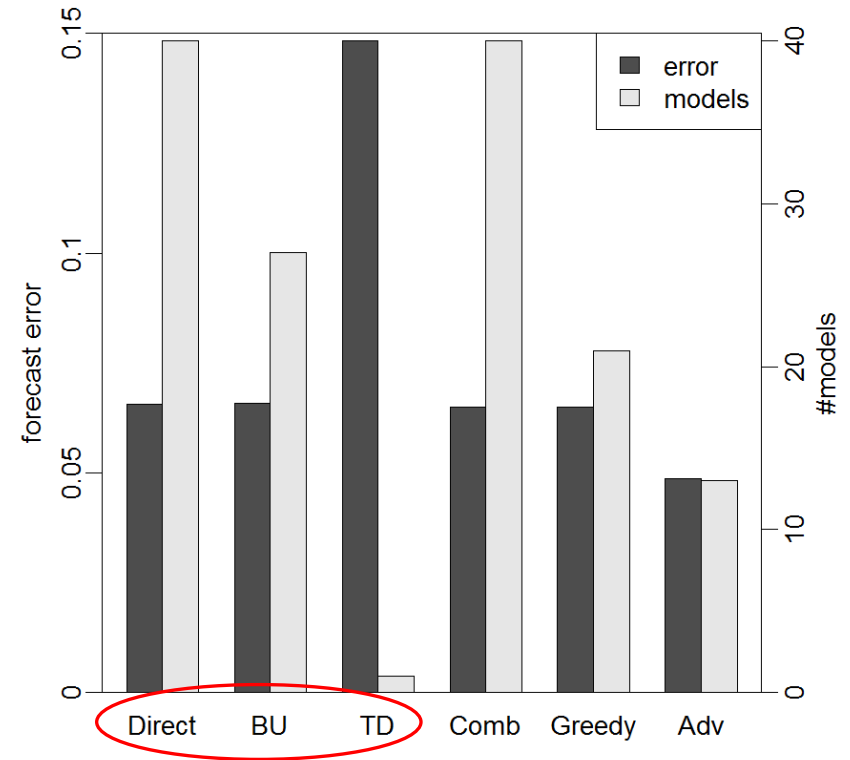
*Correlation between indicators and real forecast error*
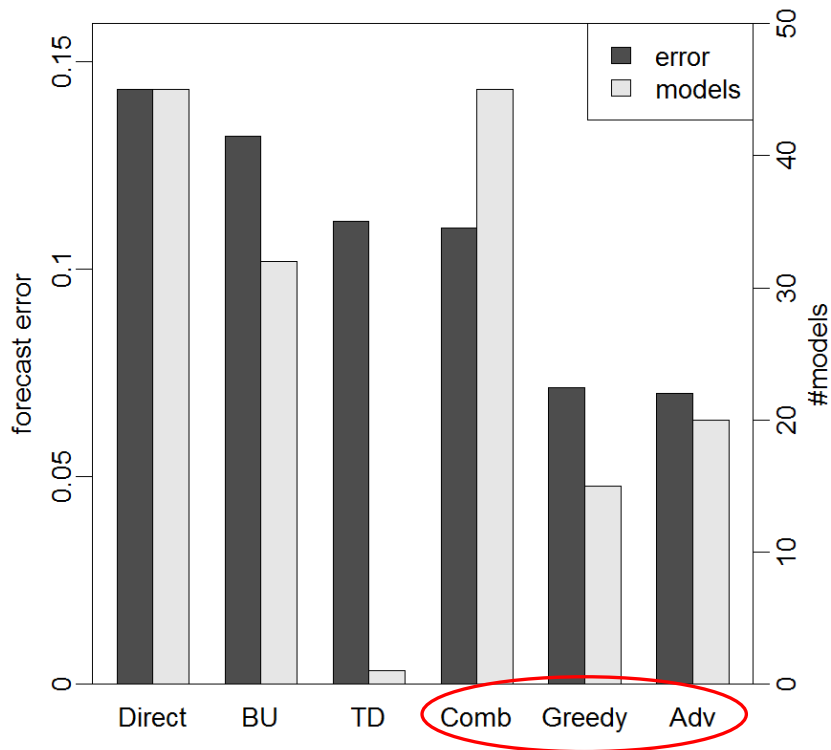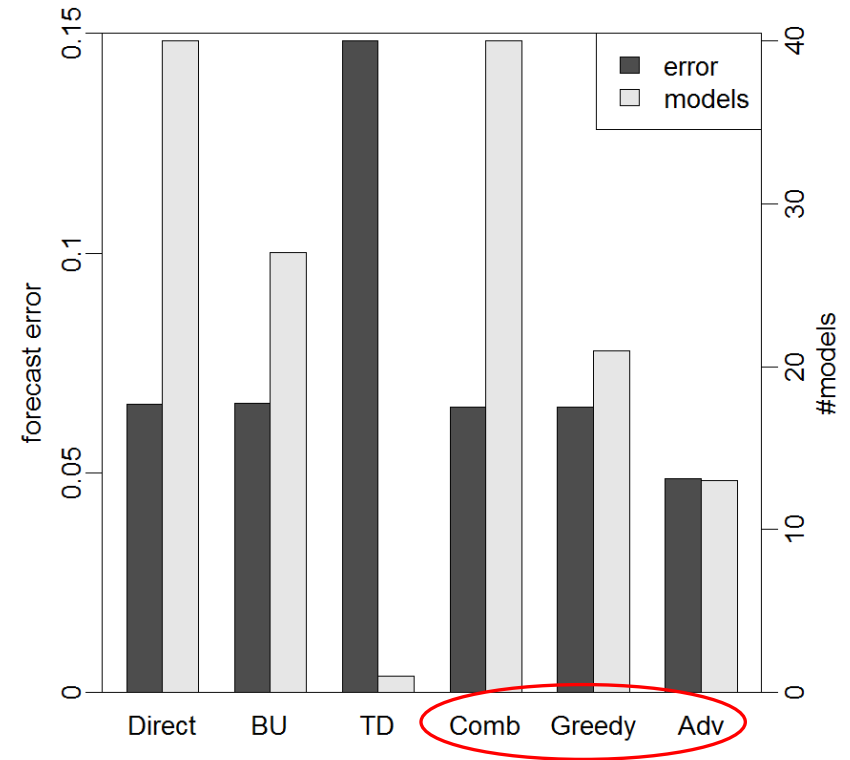
*Tourism*

*Sales*

Static approaches – data independent
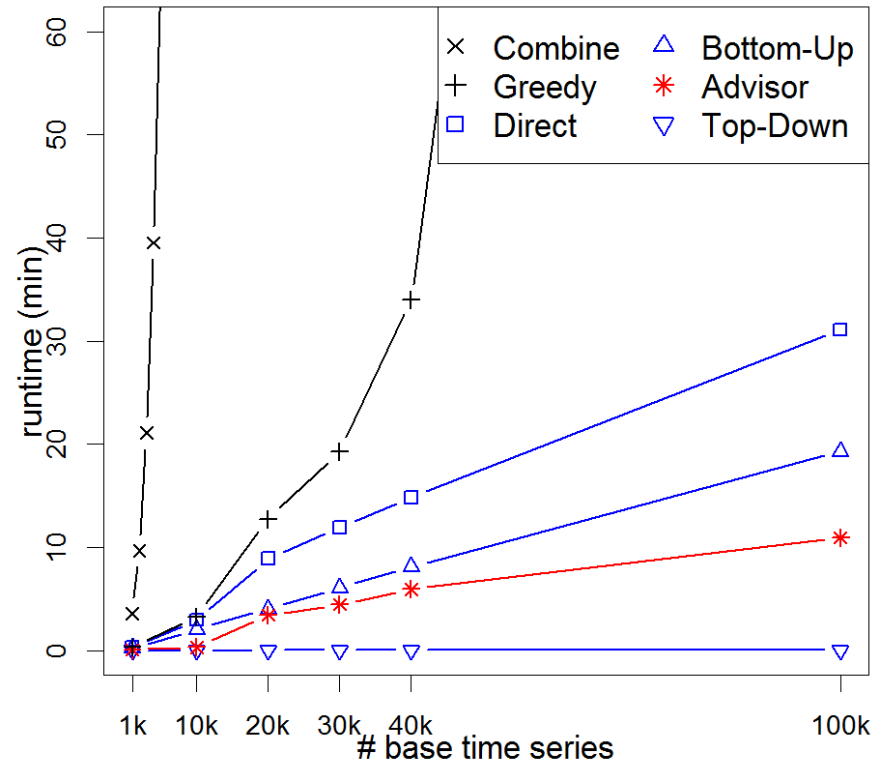
Tourism

Sales

Dynamic approaches – empirical selection

## Scalability

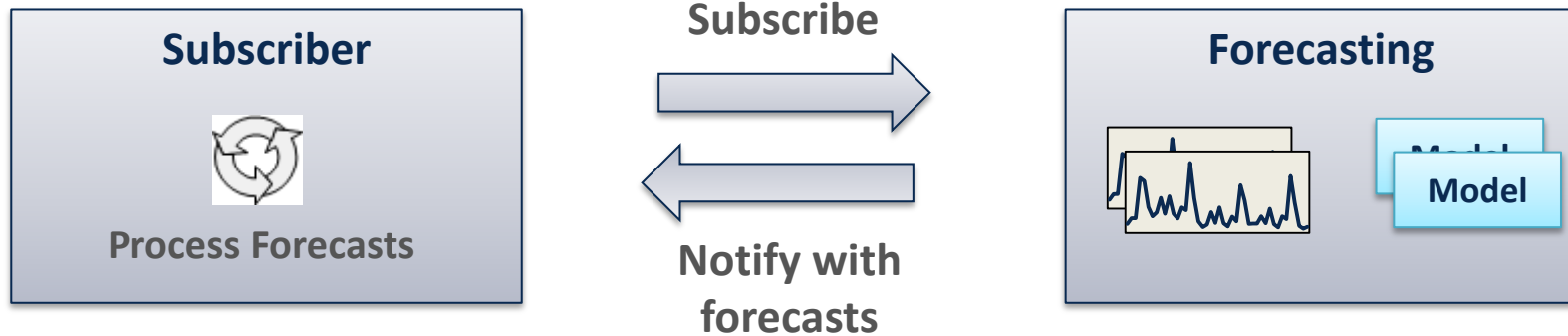# *Subscription-Based Forecast Queries*

*Energy data management systems*

- Provide stable energy supply while including larger amounts of renewable energy
- **Continously** require forecasts of energy demand and supply

**MIRABEL energy**
Balancing energy supply and demand

*Subscription-based forecast queries*

*When to notify subscriber?*

**After each new
real value …**

365, 412,
546

→ Many notifications
→ **High subscriber costs**

**As less as possible …**

365, 412,
546, 459,
460, 549,
340, 400,
128, 943,

→ Long messages
→ Low accuracy
→ Resend messages
→ **High subscriber costs**

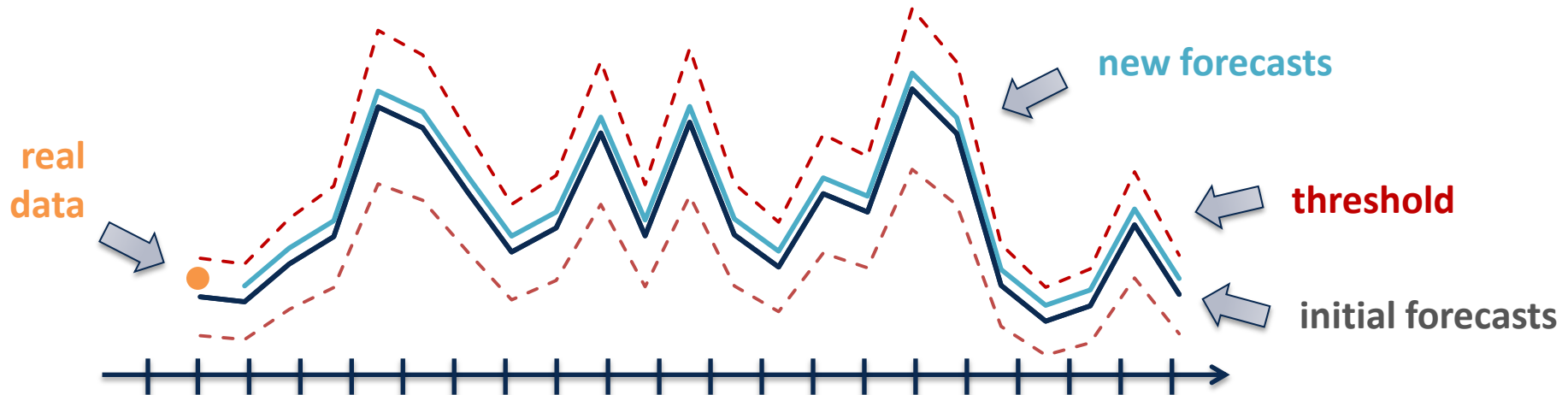**Optimal …**

365, 412,
546, 756,
349, 321

→ Reduce number of
notifications and
notification length

*Parameters*

- Time series description
- Minimum continous forecast horizon
- Accuracy threshold

```
SELECT    datetime, production
FROM      ts_powerproduction
WHERE     type = „wind"
FORECAST  3 hours
THRESHOLD 0.1
```



**new forecasts**

**real data**

**threshold**

**initial forecasts**
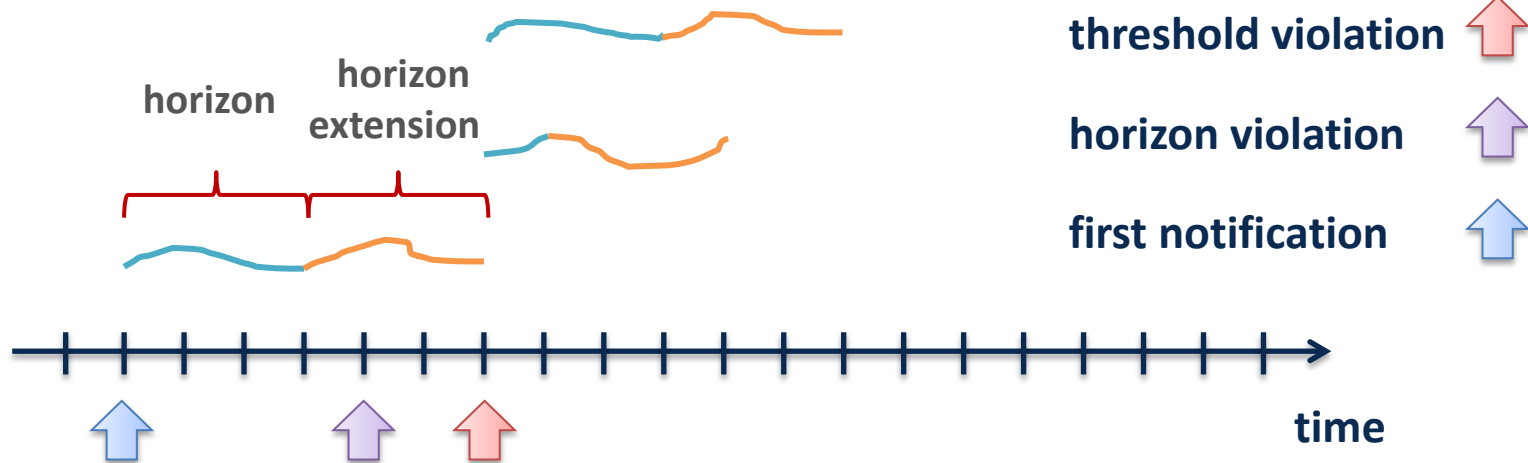
Horizon Violation

- <u>When?</u>  Subscriber has less than minimum horizon
- <u>What?</u>  Send missing values + horizon extension

Threshold Violation

- <u>When?</u>  Threshold is violated
- <u>What?</u>  Resend all values

```
SELECT      datetime, production
FROM        ts_powerproduction
WHERE       type = „wind"
FORECAST    3 hours
THRESHOLD   0.1
```
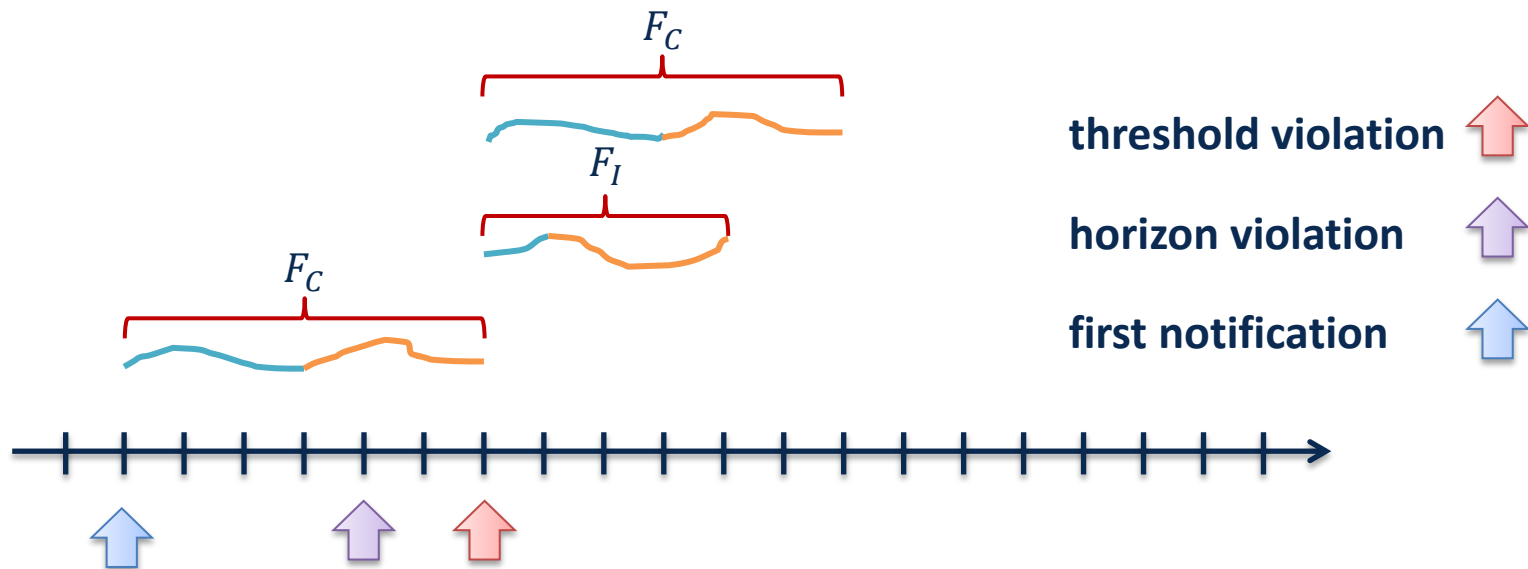
horizon

horizon
extension

threshold violation

horizon violation

first notification

time

Database **Technology**
Group

*Processing costs of the subscriber*

- Analytically known or learned function
- Depend on the forecast horizon
- Complete costs $F_C$
  - Complete restart of processing
- Incremental costs $F_I$
  - Processing of additional values

**Subscriber**

**Process Forecasts**

$F_C$

**threshold violation**

$F_I$

**horizon violation**

$F_C$

**first notification**

*Assume we know …*

- The sequence of threshold + horizon violations
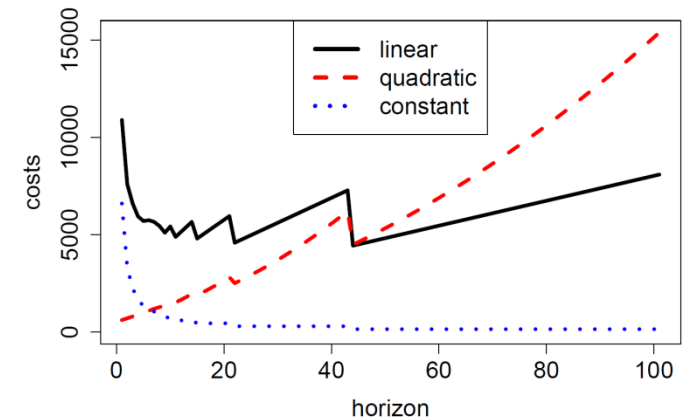- The subscriber cost functions $F_C$ and $F_I$



*Subscriber costs over subscription lifetime*

- Sum over all $F_C$ and $F_I$

*Optimization Goal*

- Find forecast horizon that minimizes total costs
- Depends on subscriber cost function
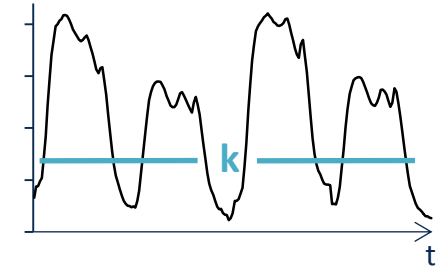- Depends on forecast accuracy



*How to get threshold violations?*

➡ Analyze past to predict future

*Core Idea: Calculate best forecast horizon using our cost model on the time series history*
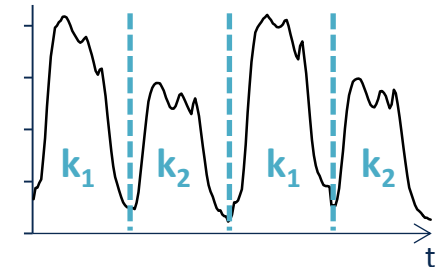
*Offline – Static*
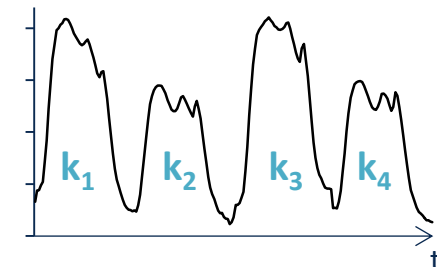
- One forecast horizon over whole lifetime

*Offline – Dynamic*

- Adapts to periodic changes of time series accuracy
- Sequence of forecast horizons for time slices

*Online*

- Adapts to arbitrary changes in data or cost functions
- Continously adapts forecast horizon

*Real-world energy demand and supply data sets*

- National Energy Demand
- Household Energy Demand
- National Wind Supply

*Subscriber cost functions*

- Synthetic linear function
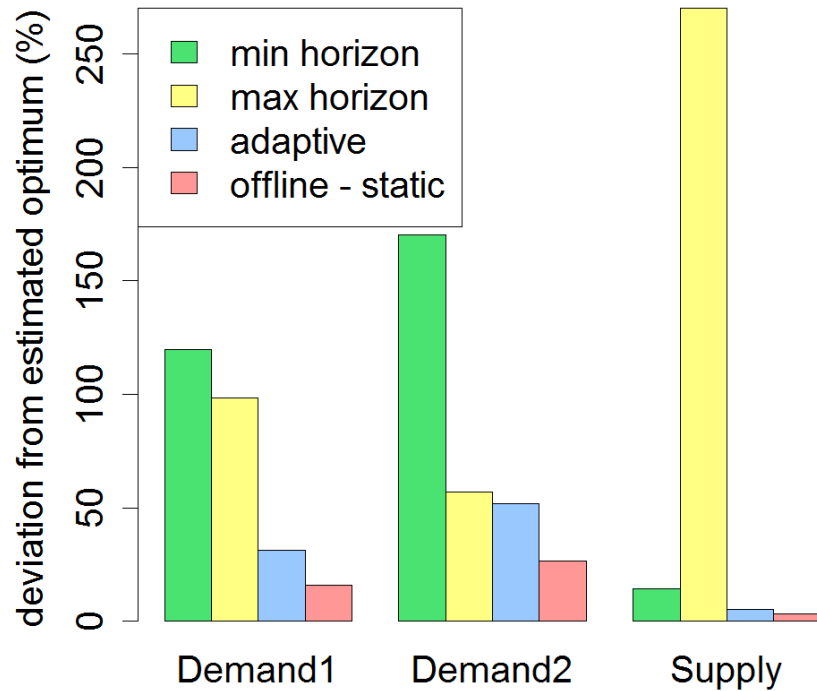- Real world cost function (obtained from MIRABEL)

*Forecast Methods*

- Tailor-made for short-term energy forecasting
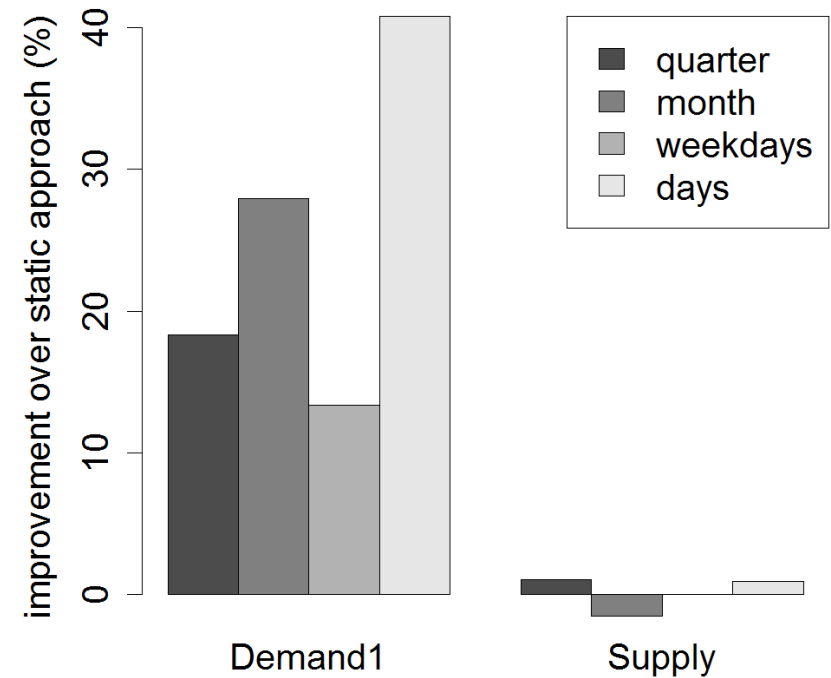- Extension of exponential smoothing

*Comparison of Approaches*

- Fixed subscription parameters and linear cost function
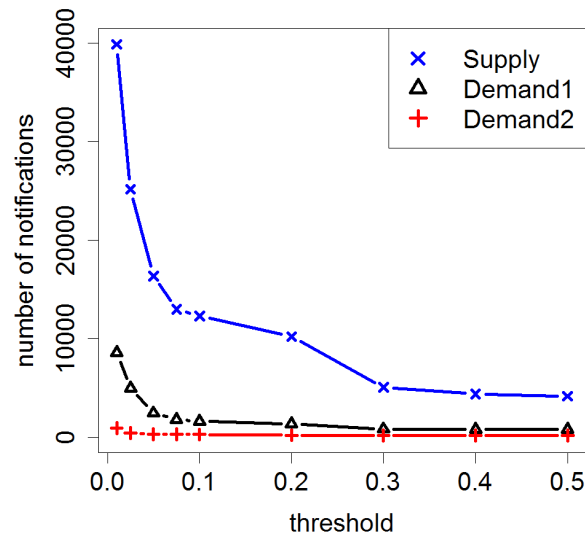
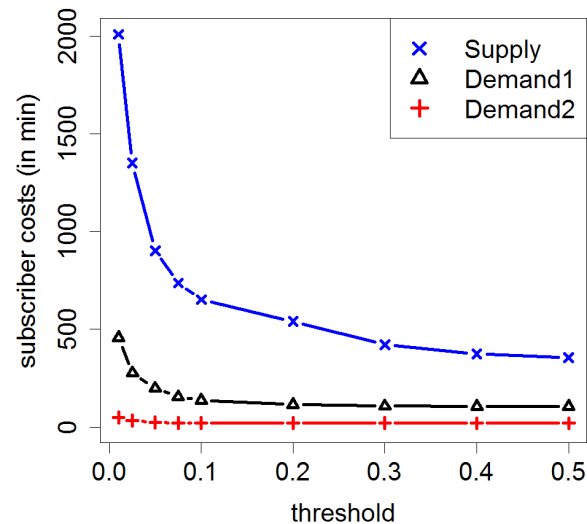*Evaluation of Time Slice Approach*

*Influence of subscriber threshold*

- Relationship between number of notifications, subscriber costs and runtime
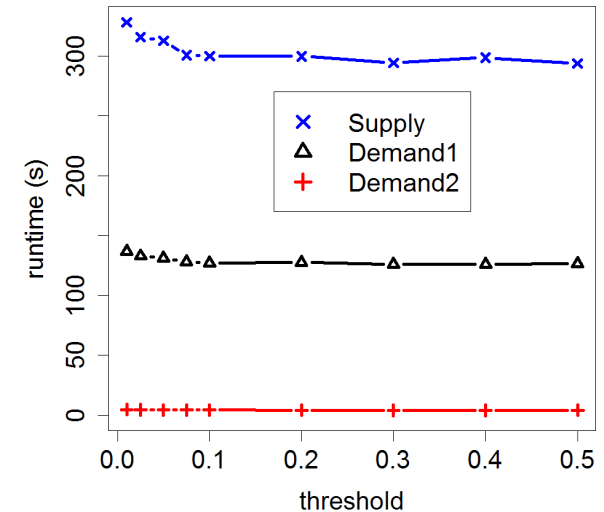- Real-world cost function

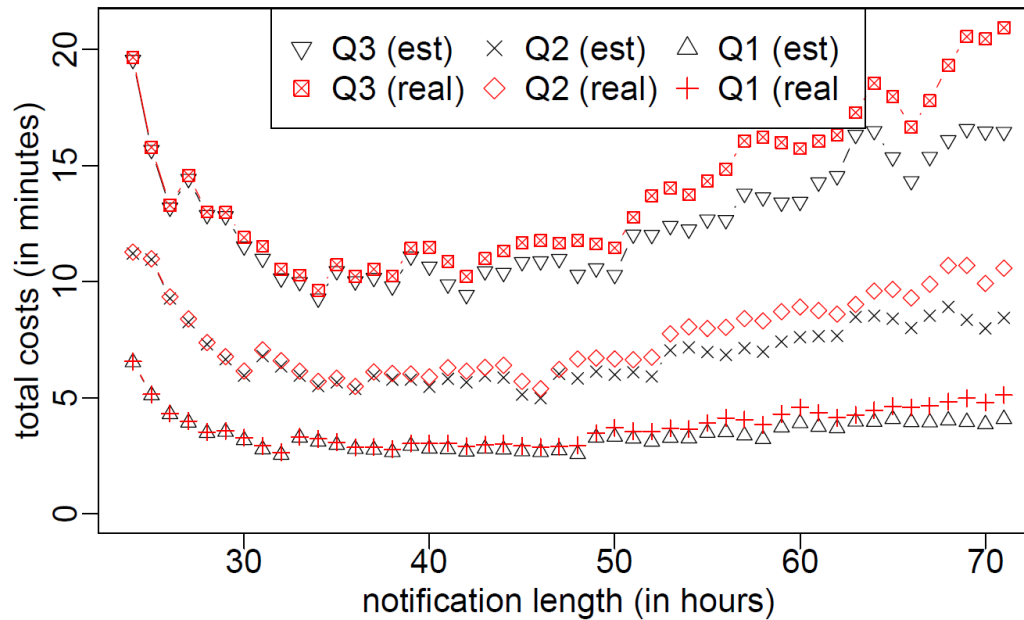number of notifications          subscriber costs                    runtime

*Cost Model Validation*

- Estimated costs vs. real costs
- Real-world cost function
- Queries with increasing complexity (Q1, Q2, Q3)

*Beyond Forecast Models*

*Towards a Model-based Database System*

## *Forecast*

- Base approach
- Predict missing data from history
- Requires no known data

## *Impute*

- Exploit local patterns
- Infer missing data from similar units
- Requires adequate set of known data
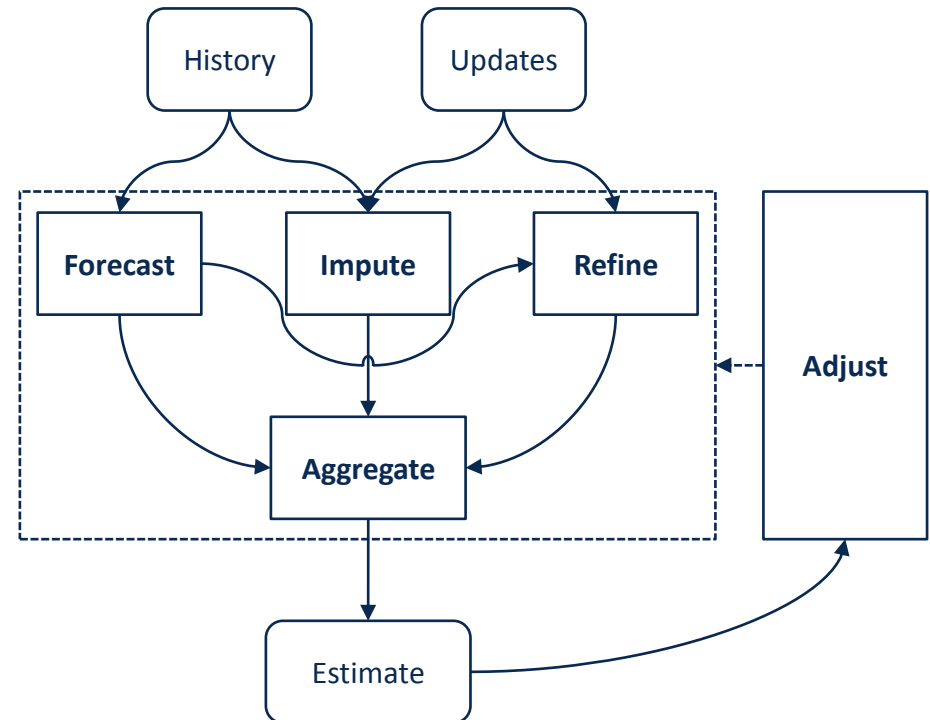
## *Refine*

- Detect local and global shifts
- Infer error → yields new forecasts
- Requires few known data

## *Aggregate*

- Calculate aggregate (e.g. report)

## *Adjust*

- Maintain models and synopses
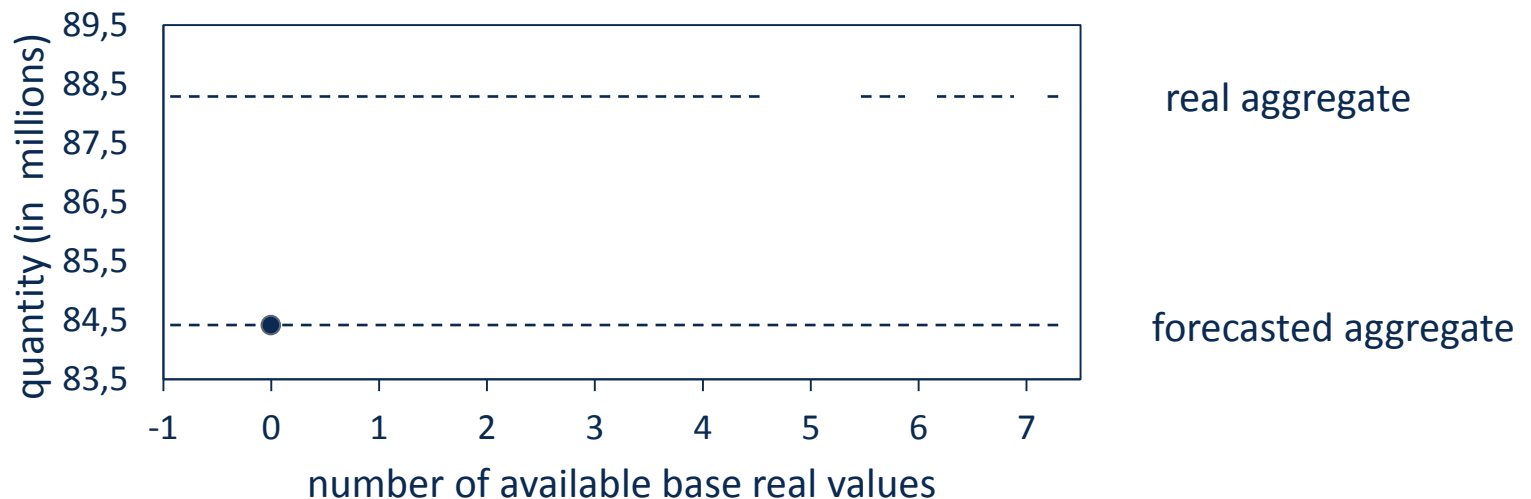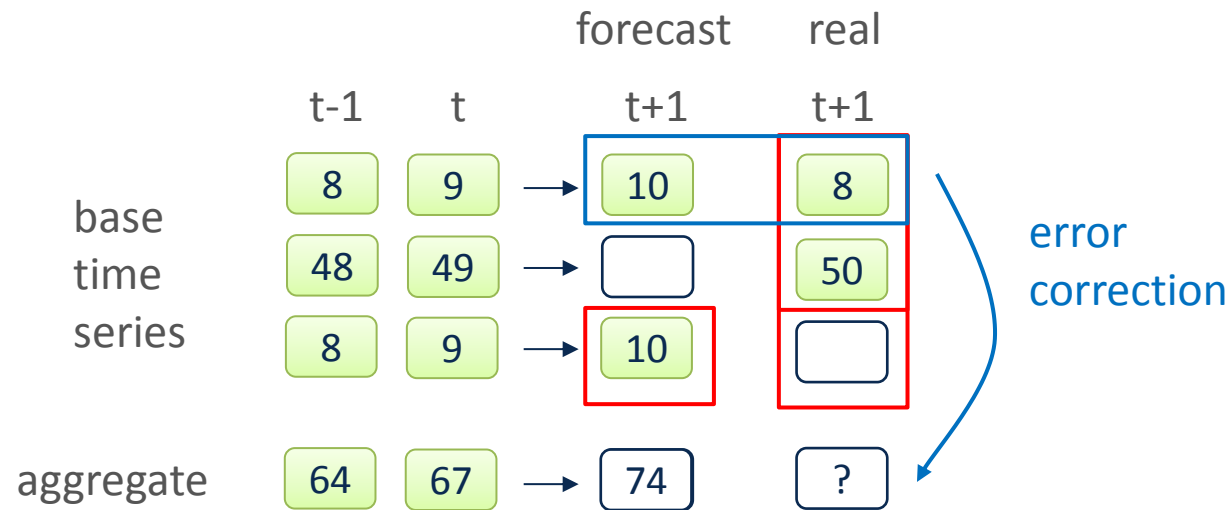- Optimize accuracy of estimates

*How to include new real data?*

*Data delivery may be late*

- There might be missing data for the last period
- Reports still have to be generated
- Estimate missing data

forecast   real

| | t-1 | t | t+1 | t+1 | |
|---|---|---|---|---|---|
| base time series | 8 | 9 | → 10 | 8 | error correction |
| | 48 | 49 | → | 50 | |
| | 8 | 9 | → 10 | | |
| aggregate | 64 | 67 | → 74 | ? | |

Refine I

$$\alpha \sum Real$$

Refine II

$$\alpha \sum (Real + Forecast)$$

Refine III

$$\alpha \sum (Real + Forecast) - \beta \sum Error$$

## *Refine III: Estimation of forecast errors*

|  | forecast | real | | error | | |
|---|---|---|---|---|---|---|
|  | t+1 | t+1 | | | | |
|  | 10 | 8 | | 8 | - | 10 |
| base | 10 | | | | ? | |
| time series | | | | | | |
|  | 50 | 50 | | 50 | - | 50 |
| aggregate | 74 | ? | | | | |

### *Case 1*

- Forecast and real value
- Calculate true model error

### *Case 2/3*

- No real value
- No error calculation

### *Case 4*

- No forecast value
- Estimate model error

*Refine – Sales*

*Refine – Wind production*

Flash Forward
Query

*FFQ project*

- Provide forecasting as 1st class citizen within a database system
- Preserve logical and physical data independence (e.g. transparent model usage, transparent model maintenance, and model creation)
- Extend traditional processing and optimization techniques
- Apply concept of traditional index advisors to foreast models

*Towards a model-based database system*

- Data is increasingly inconsistent, incomplete and imprecise
- Extend concept of models to other use cases (missing data, uncertain data, data compression …)

Wolfgang Lehner

# Forecasting and Data Imputation Strategies in Database Systems

11.07.2013