Faculty of Computer
and Information Sciences
*Ain Shams University*

كلية الحاسبات والمعلومات
جامعة عين شمس

# Data Science

**Code:**

Instructor

# Prof.Dr. Abeer M. Mahmoud

Professor of Computer Science

-faculty of Computer and Information Sciences- Ain Shams University

# Data Science and Big Data Analytics v2

# Topics : Data Science and Big Data Analytics Course

| Introduction to Big Data Analytics + Data Analytics Lifecycle | Review of Basic Data Analytic Methods Using R | Advanced Analytics – Theory and Methods | Advanced Analytics - Technology and Tools | The Endgame, or Putting it All Together + Final Lab on Big Data Analytics |
|---|---|---|---|---|
| Big Data Overview<br><br>State of the Practice in Analytics<br><br>The Data Scientist<br><br>Big Data Analytics in Industry Verticals<br><br>Data Analytics Lifecycle | Using R to Look at Data - Introduction to R<br><br>Analyzing and Exploring the Data<br><br>Statistics for Model Building and Evaluation | K-means Clustering<br><br>Association Rules<br><br>Linear Regression<br><br>Logistic Regression<br><br>Naive Bayesian Classifier<br><br>Decision Trees<br><br>Time Series Analysis<br><br>Text Analysis | Analytics for Unstructured Data (MapReduce and Hadoop)<br><br>The Hadoop Ecosystem<br><br>In-database Analytics – SQL Essentials<br><br>Advanced SQL and MADlib for In-database Analytics | Operationalizing an Analytics Project<br><br>Creating the Final Deliverables<br><br>Data Visualization Techniques<br><br>+ Final Lab – Application of the Data Analytics Lifecycle to a Big Data Analytics Challenge |

**DELL**Technologies

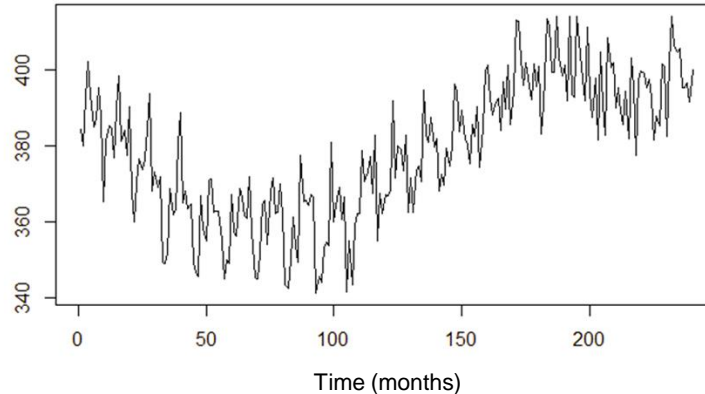# Lesson: Time series analysis

**D∕ELL**Technologies

# Time Series Analysis

During this lesson, the following topics are covered:

- Time Series Analysis and its applications in forecasting

- Autoregressive Moving Averages (ARMA) and ARIMA Models

- Implementing the Box-Jenkins Methodology using R

- Reasons to choose (+) and cautions (-) with Time Series Analysis

**D¢LL**Technologies

# Time Series Analysis, cont.

- Time Series: Ordered sequence of equally spaced values over time
- Analysis Goals
  - To identify the internal structure of the time series
  - To forecast future events
    - Example: Based on sales history, what will next December's sales be?



Time (months)
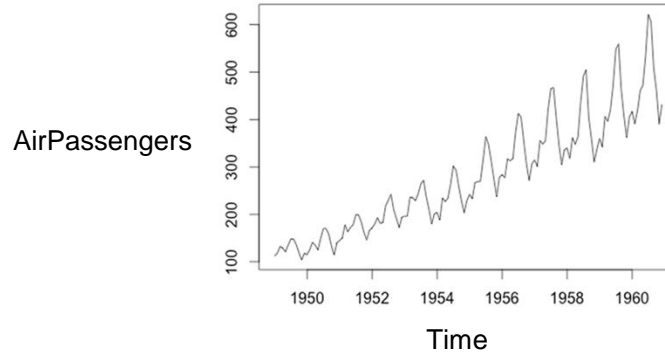
**D∢LL**Technologies

# Components of Time Series Analysis

- A time series may involve the following components:
  - Trend
  - Seasonality
  - Cycles
  - Random

- **Method: Box-Jenkins—ARMA**

**D&LL**Technologies

# Box-Jenkins method—what is it?

Models historical behavior to forecast the future

AirPassengers



Time

Applies ARMA

Accounting for the Trends and Seasonality components

Input: Time Series

Output: Expected future value of the time series

**D∕ELL**Technologies

# AR and MA models

- Autoregressive process
  - AR(p): Current value of the series depends on it own p previous values.
  - p is the order of the AR process.

- Moving Average process
  - MA(q): Current deviation from mean depends on q previous deviations.
  - q is the order of the MA process.

**D&LL**Technologies

# Time series—use cases

Forecast:

- Employment Trend

- Next month's sales

- Tomorrow's stock price

- Hourly power demand

**DELL**Technologies

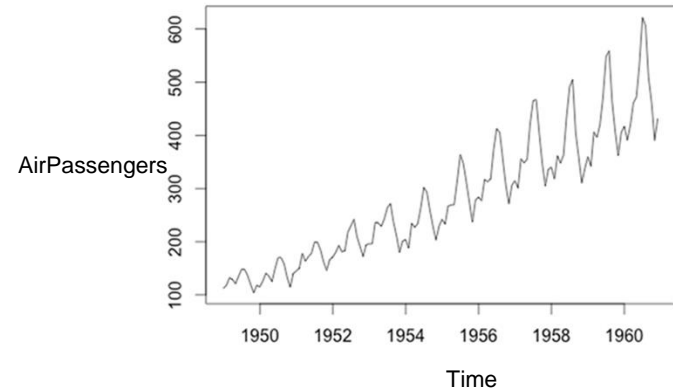# Modeling time series
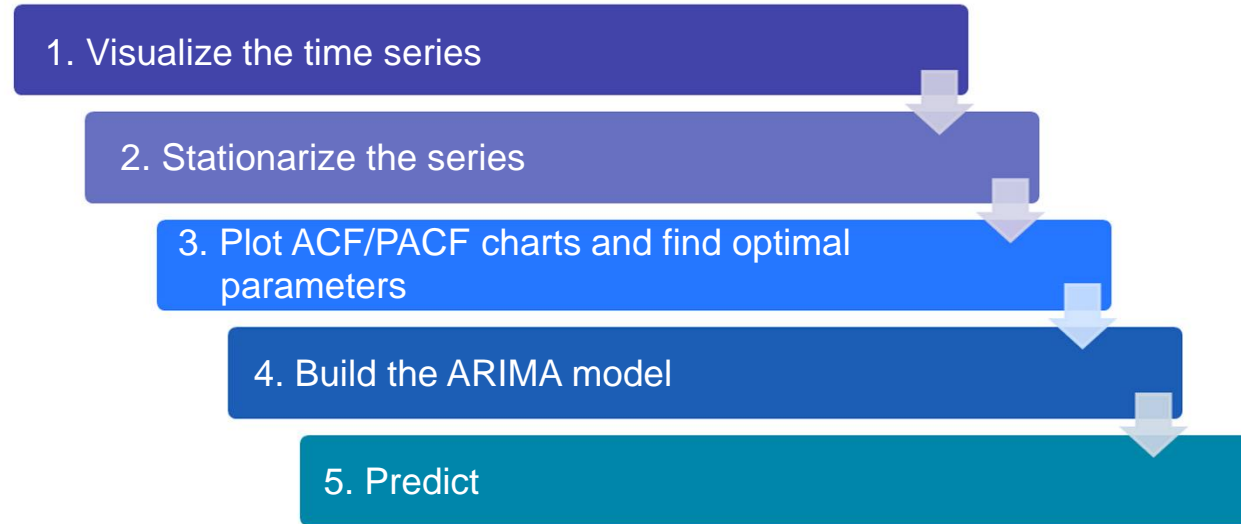
Model the time series as

$$Y_t = T_t + S_t + R_t, \qquad t=1,...,n.$$

$T_t$   Trend term
     Air travel steadily increased over the last
     few years

$S_t$   The seasonal term
     Air travel fluctuates in a regular pattern over the
     course of a year

$R_t$   Random component
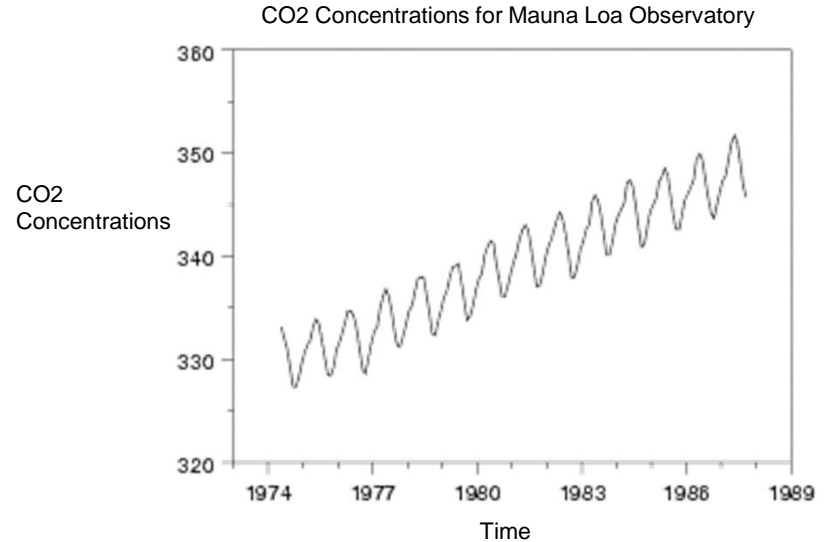     To be modeled with ARMA

AirPassengers

DELLTechnologies

# Framework for ARIMA Time Series Modeling

1. Visualize the time series

2. Stationarize the series

3. Plot ACF/PACF charts and find optimal parameters

4. Build the ARIMA model
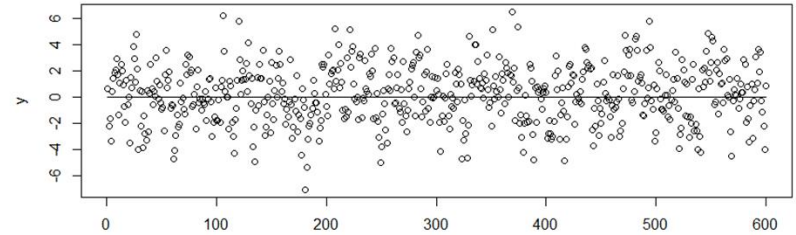
5. Predict

DELLTechnologies

# Step 1—visualizing time series

- Analyze the trend before building any time series model.

- Details of interest pertain to any patterns.
  - Trend
  - Seasonality
  - Random behavior

- This graph is an upward trending time series with seasonality every year.

CO2 Concentrations for Mauna Loa Observatory
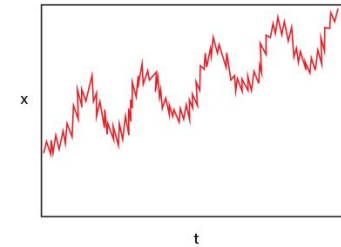
DELLTechnologies

# Step 2—stationarize series

- Box-Jenkins methodology assumes that the random component is a stationary sequence.
  - Constant mean
  - Constant variance
  - Autocorrelation does not change over time.
    - Constant correlation of a variable with itself at different times

- In practice, to obtain a stationary sequence, the data must be:
  - Detrended.
  - Seasonally adjusted.



Stationary Series



Nonstationary Series

DELLTechnologies

# Step 2—stationarize series—detrending

In this example, there is a linear trend, so fit a linear model:
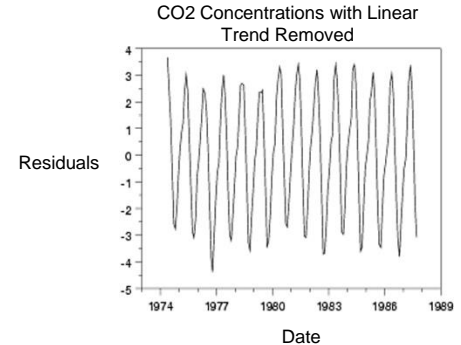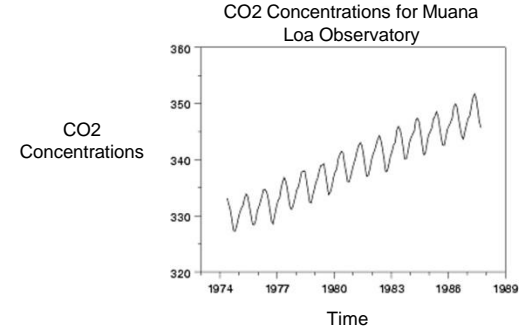
$$T_t = m \cdot t + b$$

The detrended series is then:

$$Y^1_t = Y_t - T_t$$

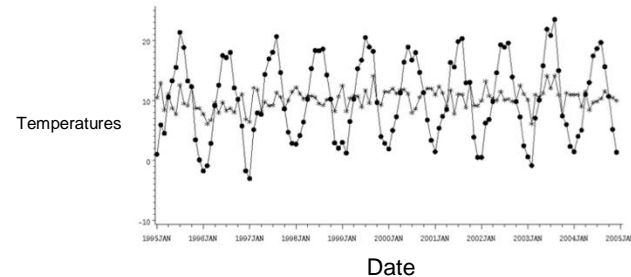Sometimes, you may have to fit a nonlinear model:

Quadratic

Exponential

CO2 Concentrations for Muana Loa Observatory

CO2 Concentrations

Time

CO2 Concentrations with Linear Trend Removed

Residuals

Date

**D⌀LL**Technologies

# Step 2—stationarize series—seasonal adjustment

Plotting the detrended series identifies seasons

> For $CO_2$ concentration, you can model the period as being a year, with variation at the month level

Simple improvised adjustment: take several years of data, calculate the average value for each month, and subtract that from $Y^1_t$

$$Y^2_t = Y^1_t - S_t$$

Box Plot of CO2 Concentrations



CO2 Concentrations

Month

Jan Feb Mar Apr May June July Aug Sep Oct Nov Dec



Temperatures

Date

DELLTechnologies

# Step 3—plot ACF and PACF to identify optimal parameters

- Auto Correlation Function (ACF)
  - Correlation of the values of the time series with itself
  - Autocorrelation "carries over"
  - Helps to determine the order, q, of a MA model
    - Where does ACF go to zero?

- Partial Auto Correlation Function (PACF)
  - An autocorrelation calculated after removing the linear dependence of the previous terms
  - Helps to determine the order, p, of an AR model
    - Where does PACF go to zero?

**DELL**Technologies

# Step 4—build model—ARMA (p, q)

$$Y_t = \boxed{\delta + \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \ldots + \phi_p Y_{t-p}}$$
$$+ \boxed{\varepsilon_t + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \ldots + \theta_q \varepsilon_{t-q}}$$

- $Y_t$ is detrended and seasonally adjusted.
- Combination of two process models
  - **Autoregressive**: $Y_t$ is a linear combination of its last *p* values.
  - **Moving average**: $Y_t$ is a constant value plus the effects of a dampened white noise process over the last *q* time values—lags.

**D∕∕LL**Technologies

# Step 4—build model—ARIMA (p, d, q)

- A stochastic modeling approach that can be used to calculate the probability of future value lying between two specified limits

- **ARIMA** adds a differencing term, $d$, to the **ARMA** model
  - Autoregressive <u>Integrated</u> Moving Average
  - Includes the detrending as part of the model:
    - Linear trend can be removed by $d$=1
    - Quadratic trend can be removed by $d$=2
    - And so on, for higher-order trends

- The general nonseasonal model is known as ARIMA ($p, d, q$):
  - $p$ is the number of autoregressive terms.
  - $d$ is the number of differences.
  - $q$ is the number of moving average terms.

**D&LL**Technologies
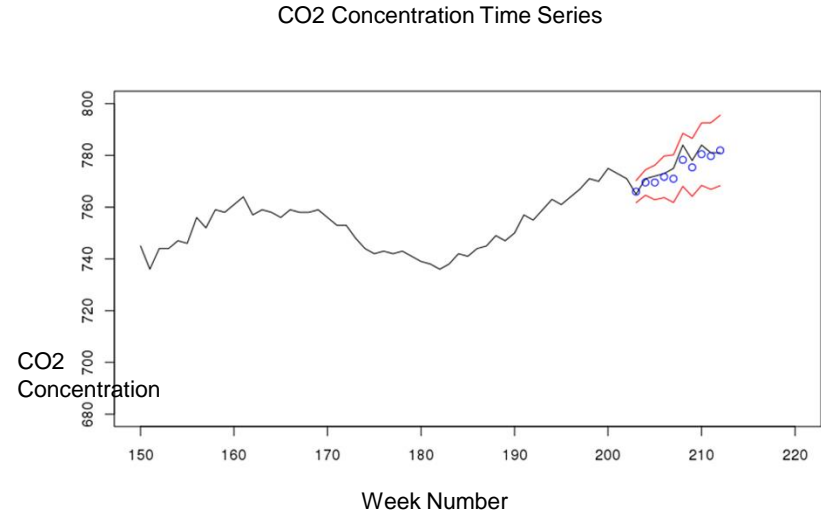
# Step 4—build model—model selection

- Based on the data, the Data Scientist selects $p$, $d$, and $q$:
  - An "art form" that requires domain knowledge, modeling experience, and a few iterations
  - Use a simple model when possible:
    - AR model ($q = 0$)
    - MA model ($p = 0$)

- Multiple models must be built and compared, using:
  - ACF and PACF.
  - p-values for the model parameter estimates
  - Information criteria metrics.
    - Imposes a penalty based on number of model parameters
    - Example: Akaike information Criterion (AIC)

**D∅LL**Technologies

# Step 5—predict

After you have the final ARIMA model, you can use this model to make future predictions with the time series.

- Here is an example output of prediction using ARIMA model. You are predicting $CO_2$ concentrations for the next four weeks, based on historical time series. The blue dots to the right of the graph indicate the prediction, and the two red lines indicate the prediction interval.



CO2 Concentration Time Series

CO2 Concentration

Week Number

**DELL**Technologies

# Time Series Analysis—reasons to choose (+) and cautions (-)

| Reasons to choose (+) | Cautions (-) |
|---|---|
| Minimal data collection<br><br>    Collect the series itself.<br><br>    It is not necessary to input drivers. | No meaningful drivers: prediction based only on past performance<br><br>    No explanatory value<br><br>    Cannot do "what-if" scenarios<br><br>    Cannot stress test |
| Designed to handle the inherent autocorrelation of lagged time series | It is an "art form" to select appropriate parameters |
| Accounts for trends and seasonality | Only suitable for short-term predictions |

**D∕ELL**Technologies

# Check your knowledge

1. What is a time series, and what are the key components of a time series?

2. How do you detrend time series data?

3. What makes data stationary?

4. How is seasonality removed from the data?

5. What are the modeling parameters in ARIMA?

6. How do you use ACF and PACF to determine the stationarity of time series data?

**D∆LL**Technologies

# Time series analysis—summary

During this lesson, the following topics were covered:

- Time series analysis and its applications in forecasting

- ARMA and ARIMA Models

- Implementing the Box-Jenkins methodology using R

- Reasons to choose (+) and cautions (-) with time series analysis

**D∕ELL**Technologies

# Module summary

| Key topics covered in this module | Methods covered in this module |
|---|---|
| Algorithms and technical foundations | Categorization, unsupervised:<br>K-means clustering<br>Association rules |
| Key use cases | Regression<br>Linear<br>Logistic |
| Diagnostics and validation of the model | Classification, supervised:<br>Naïve Bayesian classifier<br>Decision trees |
| Reasons to choose (+) and cautions (-) of the model | Time series analysis |
| Fitting, scoring, and validating model in R and in-db functions | Text analysis |

**DELL**Technologies

# Lesson: Improving model performance

**D∕LL**Technologies

# Evaluating and improving model performance

Some common situations that data scientists may encounter when it comes to model performance:

- Underfitting: The model poorly explains the training dataset

- Overfitting: The model fits the training dataset very well (almost too well), but the predictive power of the model is poor

- Once in production, the predictive power of a model degrades over time

**DELL**Technologies

# Underfitting

In the Model Building phase of an analytic project

- Underfitting is essentially the starting state
  - One or two obvious variables are added to the model
  - But the dataset is so poorly explained by the fitted model, there is no point to use the testing dataset

- To improve the model performance
  - Include additional features/variables to the model
  - Account for any correlated input variables
  - Apply necessary transformations to the input or outcome variables
  - Consider splitting one massive model into several smaller models
  - Try a different analytical model or an ensemble of analytical techniques

- Evaluate the fitted model against a testing dataset

**DELL**Technologies

# Overfitting

- Overfitting is often observed when
  - The fitted model explains the training dataset very well
  - But the model fails to accurately predict the outcomes when presented with a new dataset

- Possible causes
  - Too many unimportant variables are added to the model
  - Fine tuning (training) the model based on the testing dataset

- Possible remediation steps
  - Revisit the variable/feature selection process
  - Ensure that any unimportant variables are removed from the model
  - Look for outliers that may be unduly influencing the modeling process

**D&LL**Technologies

# Other reasons that a model may not perform well

- The datasets were not properly obtained or built
  - Proper randomization techniques need to be applied
  - oThe data is not representative of the population of interest

- Predicting outcomes for inputs outside the range of the training dataset

- The operating environment has changed since the model went into production
  - External change examples
    - New competitors enter the market
    - New product alternatives are available
    - Macroeconomic changes occur
  - Internal change examples
    - Corporate strategy changes
    - Improvement initiatives implemented
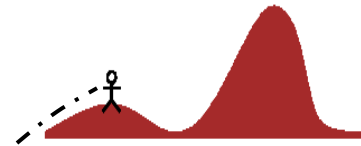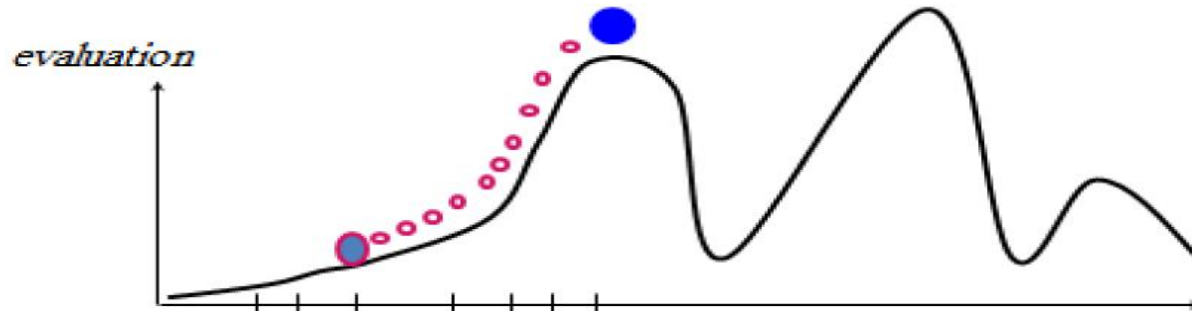    - New marketing or sales incentives

DELLTechnologies

# SUPPORT STRATEGY

- ➢ Hill-Climbing Search.
- ➢ Simulated Annealing Search.
- ➢ Local Beam Search.
- ➢ Genetic Algorithms.

31

**DELL**Technologies

# **Hill-Climbing Search**

**D∕ELL**Technologies

# Hill-Climbing Search

➢ **Main Idea**: Keep a single current node and move to a neighboring state to improve it.

➢ Uses a loop that continuously moves in the direction of increasing value (**uphill**):

➢ Choose the best successor, choose randomly if there is more than one.

➢ Terminate when a peak reached where no neighbor has a higher value.

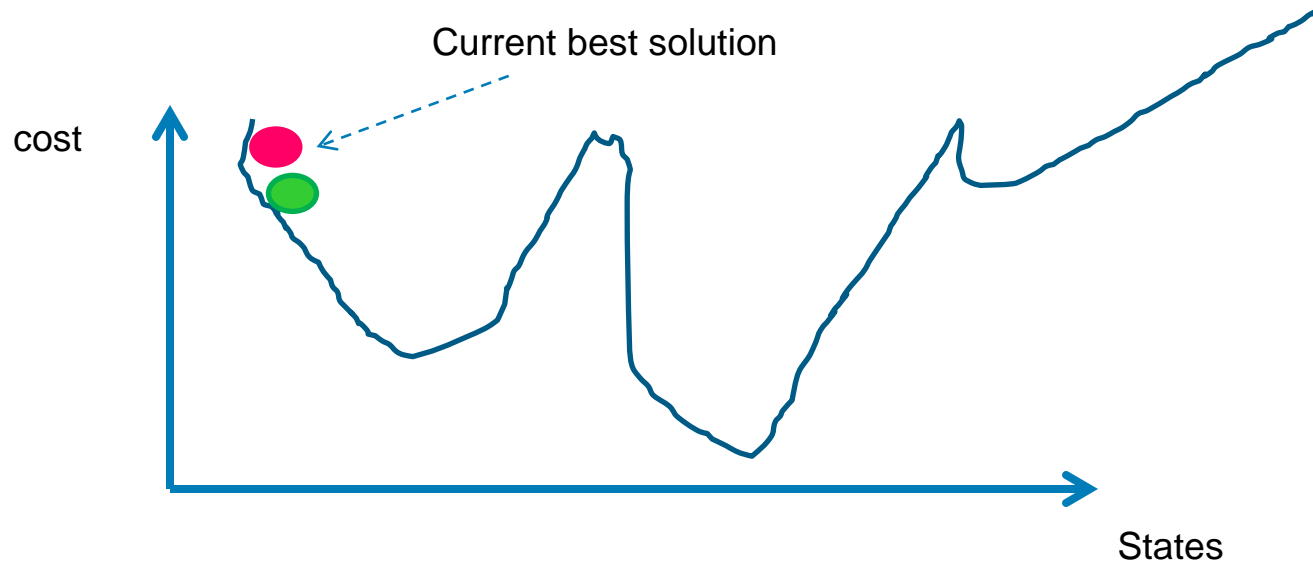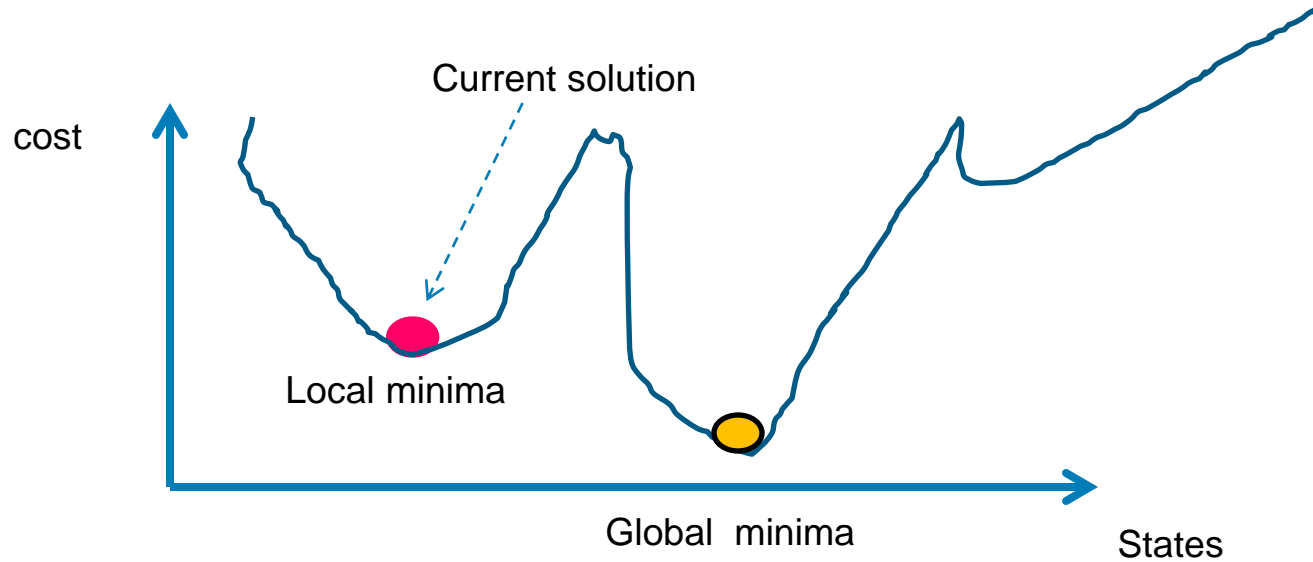➢ It also called **greedy local search**, steepest ascent/descent.

© C    33

# Hill-Climbing Search

function HILL-CLIMBING( *problem*) **returns** a state that is a local maximum
    **inputs**: *problem*, a problem
    **local variables**: *current*, a node
                      *neighbor*, a node

    $current \leftarrow$ MAKE-NODE(INITIAL-STATE[*problem*])
    **loop do**
        $neighbor \leftarrow$ a highest-valued successor of *current*
        **if** VALUE[neighbor] $\leq$ VALUE[current] **then return** STATE[*current*]
        $current \leftarrow neighbor$

34

**D∕ELL**Technologies

# Hill-Climbing in Action …



Current best solution

cost

States

35

DELLTechnologies

# Hill-Climbing in Action …



cost

Current solution

Local minima

Global minima

States

36

**D∕ELL**Technologies

# Hill-Climbing in Action …



**Drawback**: Depending on initial state, it can get stuck in local maxima/minimum or flat local maximum and not find the solution.

**Solution** :  Random restart.

37

**D∕LL**Technologies

# Simulated Annealing Search

**D∉LL**Technologies

# The Problem

❑ Most minimization strategies find the *nearest* local minimum

❑ Standard strategy

✓ Generate trial point based on current estimates

✓ Evaluate function at proposed location

✓ Accept new value if it improves solution

❑ We need a strategy to find other minima

❑ This means, we must sometimes select new points that do not improve solution   How?

**D&LL**Technologies

# Simulated annealing Search

➢ **Main Idea**: escape local maxima by allowing some "bad" moves but gradually **decrease their frequency**.

➢ Instead of picking the **best** move, it picks a **random** move..

40

**D∞LL**Technologies

# Simulated annealing Search

```
function SIMULATED-ANNEALING( problem, schedule) returns a solution state
    inputs: problem, a problem
            schedule, a mapping from time to "temperature"
    local variables: current, a node
                     next, a node
                     T, a "temperature" controlling prob. of downward steps

    current ← MAKE-NODE(INITIAL-STATE[problem])
    for t ← 1 to ∞ do
        T ← schedule[t]
        if T = 0 then return current
        next ← a randomly selected successor of current
        ΔE ← VALUE[next] − VALUE[current]
        if ΔE > 0 then current ← next
        else current ← next only with probability e^(ΔE/T)
```

Similar to hill climbing, but a random move instead of best move.

Case of improvement, make the move. Otherwise, choose the move with probability that decreases exponentially with the "badness" of the move.

➤ say the change in objective function is $\delta$
➤ **if** $\delta$ is positive, then move to that state
➤ **otherwise**:
  • move to this state with probability proportional to $\delta$
  • thus: worse moves (very large negative $\delta$) are executed less often

41

# Simulated annealing Search

– Lets say there are 3 moves available, with changes in the objective function of d1 = -0.1, d2 = 0.5,  d3 = -5. (Let T = 1).

– pick a move randomly:

  o if d2 is picked, move there.  (don't calculate just accept cause it is better)

  o if d1 or d3 are picked, probability of move = exp(d/T)

  o move 1: prob1 = exp(-0.1) = 0.9,

    ▪ i.e., 90% of the time we will accept this move

  o move 3: prob3 = exp(-5) = 0.0067

    ▪ i.e., 0.6% of the time we will accept this move

**D\u2044LL**Technologies

# Test yourself

if the T=3 , calculate the probability of move to the current state given the change in objective function in the following table

| ΔE | Returned Solution | |
|---|---|---|
| | Next | Current |
| d1 = 0.4 | _ _ _  _ _ _ _ _ | _  _ _ _ _ _ |
| d2 = - 0.2 | _ _ _ _ _ _ _ | _ _ _ _ _ _ _ |
| d3= -3 | _ _ _ _ _ _ _ _ | _ _ _ _ _ _ _ |

43

**D∕ELL**Technologies

# Properties of Simulated Annealing

- Cooling Schedule: determines rate at which the temperature T is lowered.

44

**D∕∕LL**Technologies

# Local Beam Search
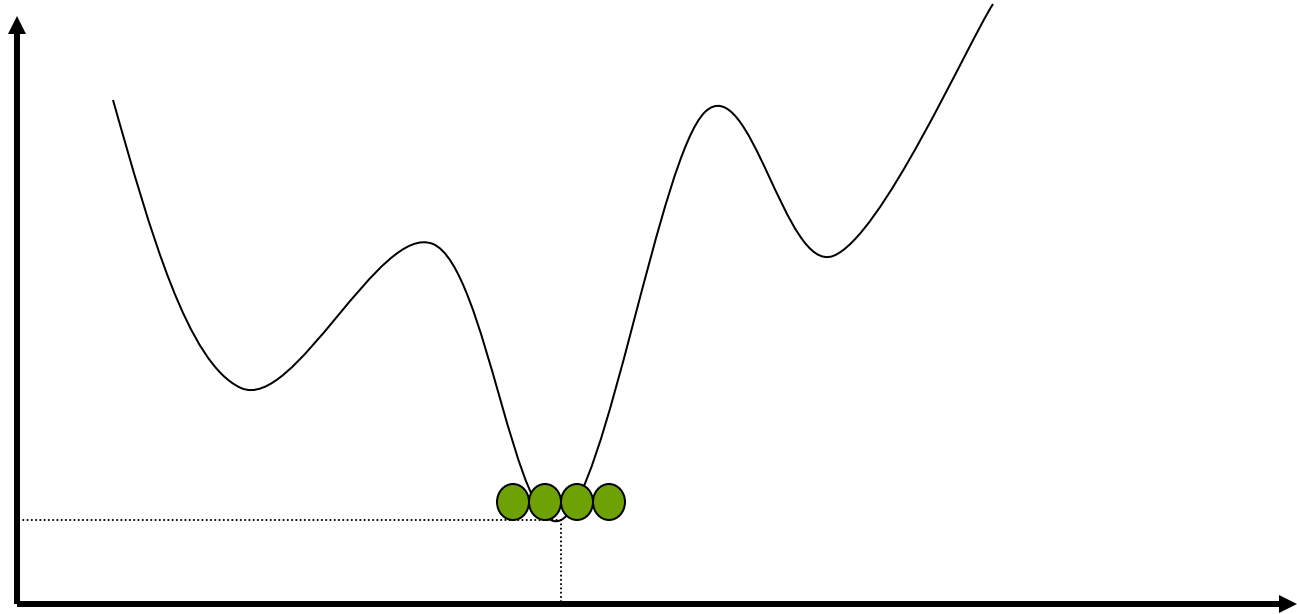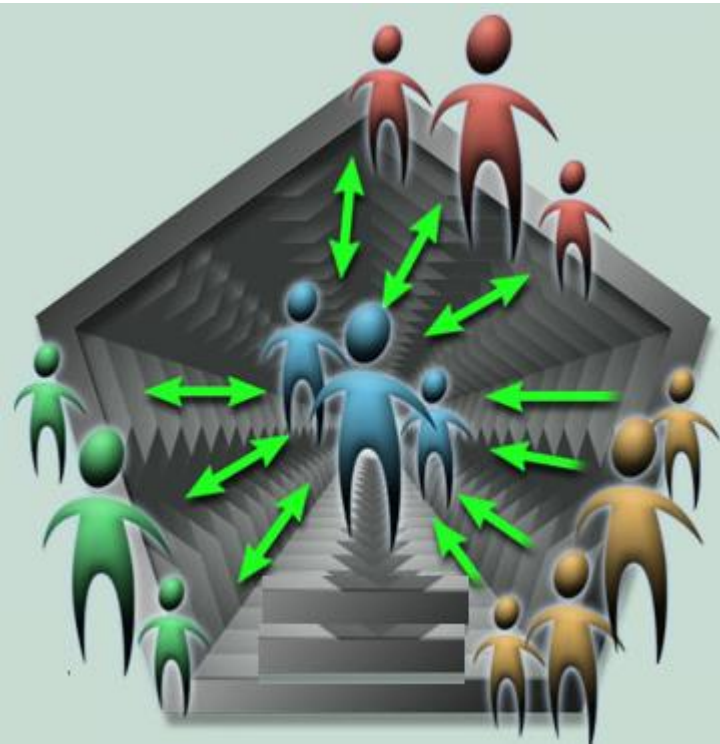
45

**D&LL**Technologies

# Local Beam Search

➢ **Main Idea**: Keep track of **k** states rather than just one.
➢ Start with **k** randomly generated states.

➢ At each iteration, all the successors of all **k** states are generated.

➢ If any one is a goal state, stop; else select the k best successors from the complete list and repeat.

➢ Drawback: the k states tend to regroup very quickly in the same region → lack of diversity.

46

DELLTechnologies

# Local Beam Search

47

DELLTechnologies

# Local Beam Search

**D&LL**Technologies

# Thank you !

**DELL**Technologies