

Statistics: A Primer

Lecture BigData Analytics

Julian M. Kunkel

julian.kunkel@gmail.com

University of Hamburg / German Climate Computing Center (DKRZ)

2017-11-27



Disclaimer: Big Data software is constantly updated, code samples may be outdated.

Outline

- 1 Descriptive Statistics
- 2 Distribution of Values
- 3 Inductive Statistics
- 4 Summary

Statistics: Overview

Statistics is the study of the collection, analysis, interpretation, presentation, and organization of data [21]

Either **describe** properties of a sample or **infer** properties of a population

Important terms [10]

- **Unit of observation:** the entity described by the data, e.g., people
- **Unit of analysis:** the major entity that is being analyzed
 - Example: Observe income of each person, analyse differences of countries
- **Statistical population:** Complete set of items that share at least one property that is subject of analysis
 - Subpopulation share additional properties, e.g., gender of people
- **Sample:** (sub)set of data collected and/or selected from a population
 - If chosen properly, they can represent the population
 - There are many sampling methods, we can never capture ALL items
- **Independence:** one observation does not effect another
 - Example: select two people living in Germany randomly
 - Dependent: select one household and pick a married couple

Statistics: Variables

- **Dependent variable:** represents the output/effect
 - Example: word count of a Wikipedia article; income of people
- **Independent variable:** assumed input/cause/explanation
 - Example: number of sentences; age, educational level

Characterization

- **Univariate analysis:** characterize a single variable
- **Bivariate analysis:** describes/analyze relationships between two vars
- **Multivariate statistics:** analyze/observe multiple dependent variables
 - Example: chemicals in the blood stream of people or chance for cancers, Independent variables are personal information/habits

Descriptive Statistics [10]

- The discipline of quantitatively describing main features of sampled data
 - Summarize observations/selected samples
- **Exploratory data analysis (EDA)**: approach for inspecting data
 - Using different chart types, e.g., box plots, histograms, scatter plot
- Methods for univariate analysis
 - Distribution of values, e.g., mean, variance, quantiles
 - Probability distribution and density
 - t-test, e.g., check if data is t-distributed
- Methods for bivariate analysis
 - Correlation coefficient¹ describes linear relationship
 - Rank correlation²: extent by which one variable increases with another var
- Methods for multivariate analysis
 - Principal component analysis (PCA) converts correlated variables into linearly uncorrelated variables called principal components

¹Pearson's product-moment coefficient

²By Spearman or Kendall

Example Dataset: Iris Flower Data Set

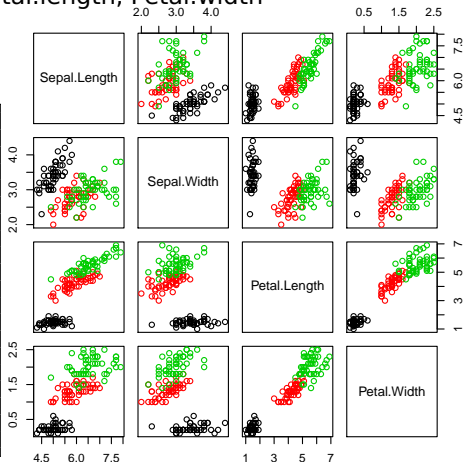
- Contains information about iris flower
- Three species: Iris Setosa, Iris Virginica, Iris Versicolor
- Data: Sepal.length, Sepal.width, Petal.length, Petal.width

R example

```

1 > data(iris) # load iris data
2 > summary(iris)
3 Sepal.Length Sepal.Width Petal.Length
4 Min. :4.300 Min. :2.000 Min. :1.000
5 1st Qu.:5.100 1st Qu.:2.800 1st Qu.:1.600
6 Median :5.800 Median :3.000 Median :4.350
7 Mean :5.843 Mean :3.057 Mean :3.758
8 3rd Qu.:6.400 3rd Qu.:3.300 3rd Qu.:5.100
9 Max. :7.900 Max. :4.400 Max. :6.900
10
11 Petal.Width Species
12 Min. :0.100 setosa :50
13 1st Qu.:0.300 versicolor:50
14 Median :1.300 virginica :50
15 Mean :1.199
16 3rd Qu.:1.800
17 Max. :2.500
18
19 # Draw a matrix of all variables
20 > plot(iris[,1:4], col=iris$Species)

```



R plot of the iris data

Distribution of Values: Histograms [10]

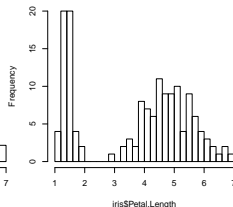
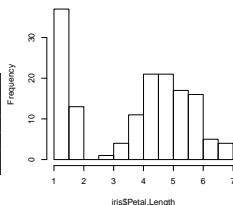
- Distribution: frequency of outcomes (values) in a sample
 - Example: Species in the Iris data set
 - setosa: 50
 - versicolor: 50
 - virginica: 50
- Histogram: graphical representation of the distribution
 - Partition observed values into bins
 - Count number of occurrences in each bin
 - It is an estimate for the probability distribution

R example

```

1 # nclass specifies the number of bins
2 # by default, hist uses equidistant bins
3 hist(iris$Petal.Length, nclass=10, main="")
4
5 hist(iris$Petal.Length, nclass=25, main="")

```



Histograms with 10 and 25 bins

Distribution of Values: Density [10]

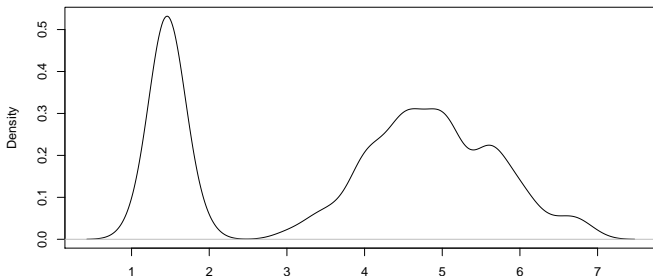
- Probability density function (density):
 - Likelihood for a **continuous** variable to take on a given value
 - Kernel density estimation (KDE) approximates the density

R example

```

1 # The kernel density estimator moves a function (kernel) in a window across samples
2 # With bw="SJ" or "nrd" it automatically determines the bandwidth, i.e., window size
3 d = density(iris$Petal.Length, bw="SJ", kernel="gaussian")
4 plot(d, main="")

```



N = 150 Bandwidth = 0.192

Density estimation of Petal.Length

Distribution of Values: Quantiles [10]

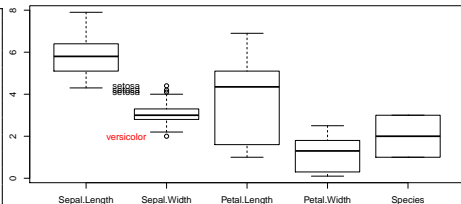
- Percentile: value below which a given percentage of observations fall
- q-Quantiles: values that partition a ranked set into q equal sized subsets
- **Quartiles**: 3 data points that split a ranked set into 4 equal points
 - $Q1=P(25)$, $Q2=\text{median}=P(50)$, $Q3=P(75)$, **interquartile range** $\text{iqr}=Q3-Q1$
- **Five-number summary**: (min, Q1, Q2, Q3, max)
- Boxplot: shows quartiles (Q1,Q2,Q3) and whiskers
 - Whiskers extend to values up to 1.5 iqr from Q1 and Q3
 - Outliers are outside of whiskers

R example

```

1 > boxplot(iris, range=1.5) # 1.5 interquartile range
2 > d = iris$Sepal.Width
3 > quantile(d)
4 0% 25% 50% 75% 100%
5 2.0 2.8 3.0 3.3 4.4
6 > q3 = quantile(d,0.75) # pick value below which are 75%
7 > q1 = quantile(d,0.25)
8 > irq = (q3 - q1)
9 # identify all outliers based on the interquartile range
10 > mask = d < (q1 - 1.5*irq) | d > (q3 + 1.5*irq)
11 # pick outlier selection from full data set
12 > o = iris[mask,]
13 # draw the species name of the outliers on the boxplot
14 > text(rep(1.5,nrow(o)), o$Sepal.Width, o$Species,
      ↪ col=as.numeric(o$Species))

```



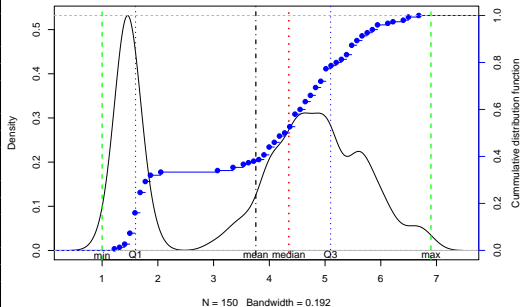
Boxplot

Density Plot Including Summary

```

1 d = density(iris$Petal.Length, bw="SJ",
  ↪ kernel="gaussian")
2
3 # add space for two axes
4 par(mar=c(5, 4, 4, 6) + 0.1)
5 plot(d, main="")
6 # draw lines for Q1, Q2, Q3
7 q = quantile(iris$Petal.Length)
8 q = c(q, mean(iris$Petal.Length))
9 abline(v=q[1], lty=2, col="green", lwd=2)
10 abline(v=q[2], lty=3, col="blue", lwd=2)
11 abline(v=q[3], lty=3, col="red", lwd=3)
12 abline(v=q[4], lty=3, col="blue", lwd=2)
13 abline(v=q[5], lty=2, col="green", lwd=2)
14 abline(v=q[6], lty=4, col="black", lwd=2)
15 # Add titles
16 text(q, rep(-0.01, 5), c("min", "Q1", "median",
  ↪ "Q3", "max", "mean"))
17
18 # identify x limits
19 xlim = par("usr")[1:2]
20 par(new=TRUE)
21
22 # Empirical cumulative distribution function
23 e = ecdf(iris$Petal.Length)
24 plot(e, col="blue", axes=FALSE, xlim=xlim, ylab="",
  ↪ xlab="", main="")
25
26 axis(4, ylim=c(0,1.0), col="blue")
27 mtext("Cumulative distribution function", side=4,
  ↪ line=2.5)

```



Density estimation with 5-number summary and cumulative density function

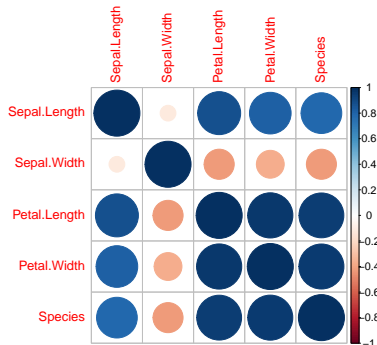
Correlation Coefficients

- Measures (linear) correlation between two variables
 - Result is between -1 and +1
 - >0.7: strong positive correlation
 - >0.2: weak positive correlation
 - 0: no correlation, < 0: negative correlation

R example

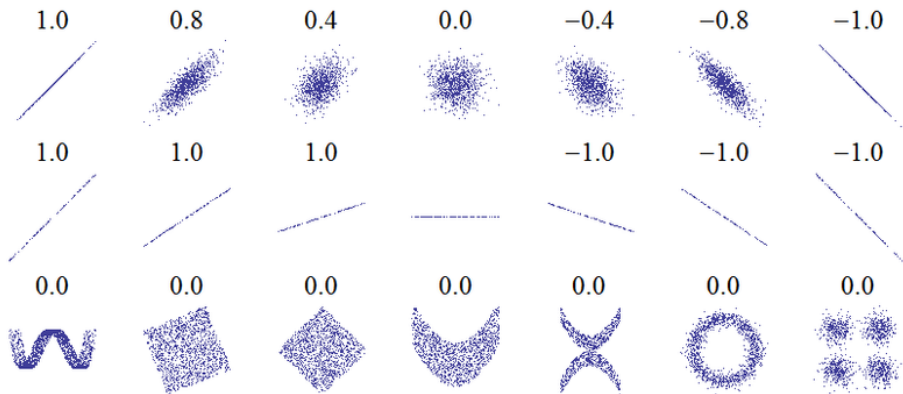
```

1 library(corrplot)
2 d = iris
3 d$Species = as.numeric(d$Species) # It is normally not
   ↪ advised to convert categorical data to numeric
4 corrplot(cor(d), method = "circle") # the right plot
5
6 mplot = function(x,y, name){
7   pdf(name,width=5,height=5) # plot into a PDF
8   p = cor(x,y, method="pearson") # compute correlation
9   k = cor(x,y, method="spearman")
10  plot(x,y, xlab=sprintf("x\n cor. coeff: %.2f rank coef.:
   ↪ %.2f", p, k))
11  dev.off()
12 }
13
14 mplot(iris$Petal.Length, iris$Petal.Width, "iris-corr.pdf")
15 # cor. coeff: 0.96 rank coef.: 0.94
16
17 x = 1:10; y = c(1,3,2,5,4,7,6,9,8,10)
18 mplot(x,y, "linear.pdf") # cor. coeff: 0.95 rank coef.: 0.95
19
20 mplot(x, x*x*x, "x3.pdf") # cor. coeff: 0.93 rank coef.: 1
  
```



Corrplot

Example Correlations for Plots of Two Variables (X, Y)



Correlations for x,y plots; Source: [22]

1 Descriptive Statistics

2 Distribution of Values

3 Inductive Statistics

4 Summary

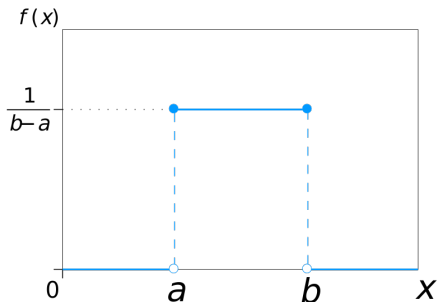
Distribution of Values

- Data exploration helps to understand distribution of values for a variable
- To illustrate data distribution, statistics uses:
 - Probability density function (PDF): value and probability; $P(x)$
 - Cumulative distribution function (CDF): sum from $-\infty$ to value; $P(X \leq x)$
- Many standard distributions exist that can be models for processes
 - Discrete probability functions are drawn from a limited set of values
 - Continuous probability functions are $\in \mathbb{R}$
- Discrete probability functions are derived from useful processes

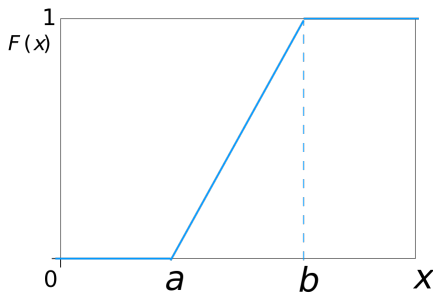
When we explore data, we may identify these common distributions

(Continuous) Uniform Distribution

- Notation: $\text{unif}(a, b)$
- Mean: $\frac{a+b}{2}$
- Variance: $\frac{(b-a)^2}{12}$



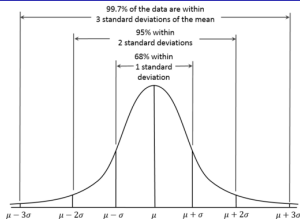
(a) Source: IkamusumeFan, PDF of the uniform probability distribution (10)



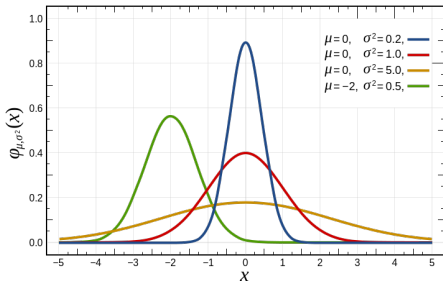
(b) Source: IkamusumeFan, CDF of the uniform probability distribution (10)

Normal Distribution

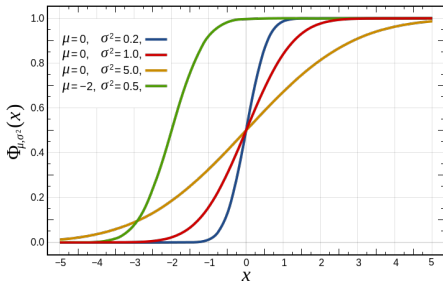
- Notation: $\mathcal{N}(\mu, \sigma^2)$
- Mean: μ
- Variance: σ^2



Source: Dan Kernler, For the normal distribution, the values less than one standard deviation away from the mean account for 68.27% of the set; while two standard deviations from the mean account for 95.45%; and three standard deviations account for 99.73%. [10]



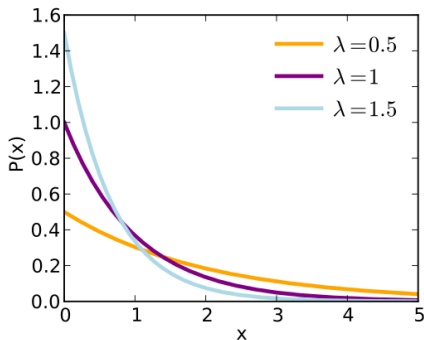
(a) Source: InductiveLoad, Probability density function for the normal distribution (10)



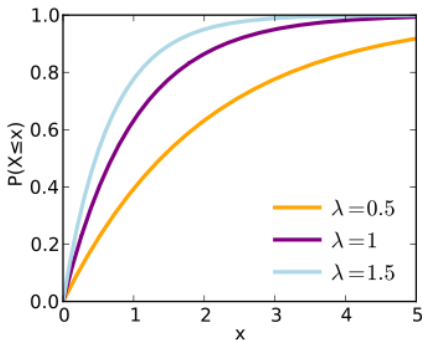
(b) Source: InductiveLoad, Cumulative distribution function for the normal distribution (10)

Exponential Distribution

- Notation: Exponential(λ)
 - λ (rate) > 0
- Mean: λ^{-1}
- Variance: λ^{-2}
- Models inter-arrival time of a Poisson process
- Memoryless, useful to model failure rates



(c) Source: Skbkakas, Probability density function (10)



(d) Source: Skbkakas, Cumulative distribution (10)

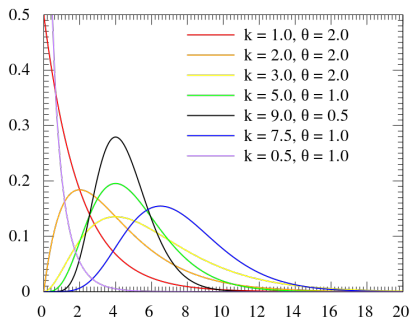
Gamma Distribution

■ Notation: $\text{Gamma}(k, \Theta)$

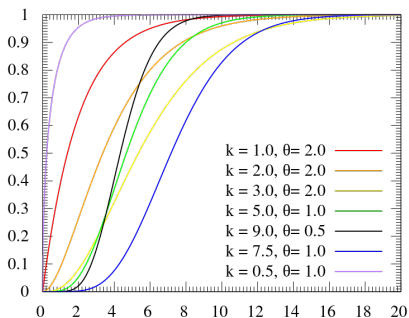
- k (shape), Θ (scale) > 0 ; sometimes rate $b = \frac{1}{\Theta}$

■ Mean: $k \cdot \Theta$

■ Variance: $k \cdot \Theta^2$



(e) MarkSweep and Cburnett, Probability density plots of gamma distributions (10)

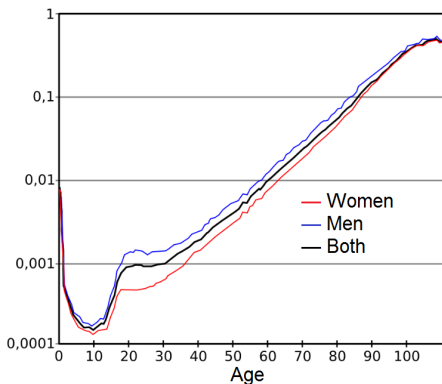


(f) MarkSweep and Cburnett, Cumulative distribution plots of gamma distributions (10)

Real Data

- Observations can be based on several (unknown) processes

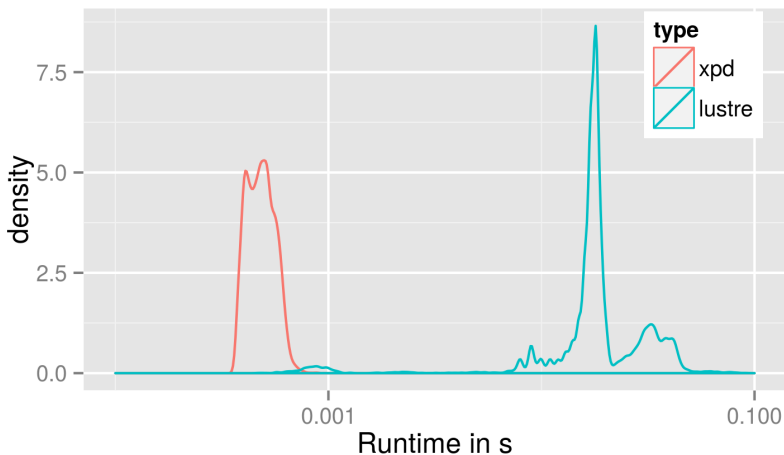
Fraction of people of a given age that died in 1999



Mark Bailin, Mortality by age [10]

Real Data (2)

Performance of reading 1 MiB of data from two different storage systems



Lustre parallel file systems vs. pooled memory of the XPD

1 Descriptive Statistics

2 Distribution of Values

3 Inductive Statistics

4 Summary

Inductive Statistics: Some Terminology [10]

- Statistical inference is the process of **deducting properties** of a population by analyzing samples
 - Build a statistical model and test the hypothesis if it applies
 - Allows to deduct propositions (statements about data properties)
- **Statistical hypothesis**: hypothesis that is testable on a process modeled via a set of random variables
- **Statistical model**: embodies a set of assumptions concerning the generation of the observed data, and similar data from a larger population. A model represents, often in considerably idealized form, the data-generating process
- **Validation**: process to verify that a model/hypothesis is likely to represent the observation/population
- **Significance**: a significant finding is one that is determined (statistically) to be very unlikely to happen by chance
- **Residual**: difference of observation and estimated/predicted value

Statistics: Inductive Statistics [10]

Testing process

- 1 Formulate default (null³) and alternative hypothesis
- 2 Formulate statistical assumptions, e.g., independence of variables
- 3 Decide which statistical tests can be applied to disprove null hypothesis
- 4 Choose significance level α for wrongly rejecting null hypothesis
- 5 Compute test results, especially the **p-value**⁴
- 6 If p-value $< \alpha$, then reject null hypothesis and go for alternative
 - Be careful: (p-value $\geq \alpha$) \nRightarrow null hypothesis is true, though it may be

Example hypotheses

- Petal.Width of each iris flowers species follow a normal distribution
- Waiting time of a supermarket checkout queue is gamma distributed

³We try to reject/**nullify** this hypothesis.

⁴Probability of obtaining a result equal or more extreme than observed.

Checking if Petal.Width is Normal Distributed

R example

```

1 # The Shapiro-Wilk-Test allows for testing if a population represented by a sample is normal distributed
2 # The Null-hypothesis claims that data is normal distributed
3
4 # Let us check for the full population
5 > shapiro.test(iris$Petal.Width)
6 # W = 0.9018, p-value = 1.68e-08
7 # Value is almost 0, thus reject null hypothesis => in the full population, Petal.Width is not normal distributed
8
9 # Maybe the Petal.Width is normal distributed for individual species?
10 for (spec in levels(iris$Species)){
11   print(spec)
12   y = iris[iris$Species==spec,]
13
14   # Shapiro-Wilk-test checks if data is normal distributed
15   print(shapiro.test(y$Petal.Width))
16 }
17
18 [1] "virginica"
19 W = 0.9598, p-value = 0.08695
20 # Small p-value means a low chance this happens, here about 8.7%
21 # With the typical significance level of 0.05 Petal.Width is normal distributed
22 # For simplicity, we may now assume Petal.Width is normal distributed for this species
23
24 [1] "setosa"
25 W = 0.7998, p-value = 8.659e-07 # it is not normal distributed
26
27 [1] "versicolor"
28 W = 0.9476, p-value = 0.02728 # still too unlikely to be normal distributed

```


Linear Models (for Regression) [10]

- **Linear regression:** Modeling the relationship between dependent var Y and explanatory variables X_i
- Assume n samples are observed with their values in the tuples $(Y_i, X_{i1}, \dots, X_{ip})$
 - Y_i is the dependent variable (label)
 - X_{ij} are independent variables
 - Assumption for linear models: normal distributed variables
- A linear regression model fits $Y_i = c_0 + c_1 \cdot f_1(X_{i1}) + \dots + c_p \cdot f_p(X_{ip}) + \epsilon_i$
 - Determine coefficients c_0 to c_p to minimize the error term ϵ
 - The functions f_i can be non-linear

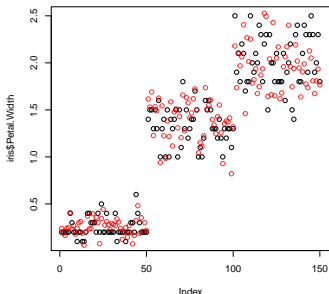
R example

```

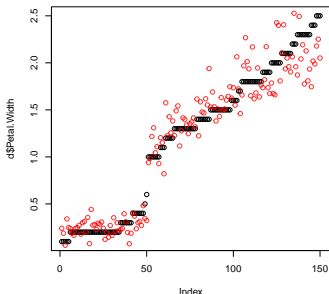
1 # R allows to define equations, here Petal.Width is our dependent var
2 m = lm("Petal.Width ~ Petal.Length + Sepal.Width", data=iris)
3
4 print(m) # print coefficients
5 # (Intercept) Petal.Length Sepal.Width
6 # -0.7065 0.4263 0.0994
7 # So Petal.Width = -0.7065 + 0.4263 * Petal.Length + 0.0994 * Sepal.Width

```

Compare Prediction with Observation



Iris linear model



With sorted data

```
1 # Predict petal.width for a given petal.length and sepal.width
2 d = predict(m, iris)
3
4 # Add prediction to our data frame
5 iris$prediction = d
6
7 # Plot the differences
8 plot(iris$Petal.Width, col="black")
9 points(iris$prediction, col=rgb(1,0,0,alpha=0.8))
10
11 # Sort observations
12 d = iris[sort(iris$Petal.Width, index.return=TRUE)$ix,]
13 plot(d$Petal.Width, col="black")
14 points(d$prediction, col=rgb(1,0,0,alpha=0.8))
```

Analysing Model Accuracy [23]

- Std. error of the estimate: variability of c_i ; should be lower than c_i
- t-value: measures how useful a variable is for the model
- $Pr(> |t|)$ two-sided p-value: probability that the variable is not significant
- Degrees of freedom: number of independent samples (avoid overfitting!)
- R-squared: fraction of variance explained by the model, 1 is optimal
- F-statistic: the f-test analyses the model goodness – high value is good

```

1 summary(m) # Provide detailed information about the model
2 # Residuals:
3 #      Min       1Q   Median       3Q      Max
4 # -0.53907 -0.11443 -0.01447  0.12168  0.65419
5 #
6 # Coefficients:
7 #              Estimate Std. Error t value Pr(>|t|)
8 # (Intercept)  -0.70648   0.15133  -4.668  6.78e-06 ***
9 # Petal.Length  0.42627   0.01045  40.804  < 2e-16 ***
10 # Sepal.Width   0.09940   0.04231   2.349   0.0201 *
11 # ---
12 # Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 " " 1
13 #
14 # Residual standard error: 0.2034 on 147 degrees of freedom
15 # Multiple R-squared:  0.9297, Adjusted R-squared:  0.9288
16 # F-statistic: 972.7 on 2 and 147 DF,  p-value: < 2.2e-16

```

- Akaike's Information Criterion (AIC)
- Idea: prefer accurate models with smaller number of parameters
- Test various models to reduce AIC
- Improve good candidates
- AIC allows to check which models can be excluded

Time Series

- A time series is a sequence of observations
 - e.g., temperature, or stock price over time
 - Prediction of the future behavior is of high interest
- An observation may depend on any previous observation
 - Trend: tendency in the data
 - Seasonality: periodic variation

Prediction models

- Autoregressive models: AR(p)
 - Depend linearly on last p values (+ white noise)
- Moving average models: MA(q)
 - Random shocks: Depend linearly on last q white noise terms
- Autoregressive moving average (ARMA) models
 - Combine AR and MA models
- Autoregressive integrated moving average: ARIMA(p, d, q)
 - Combines AR, MA and differencing (seasonal) models

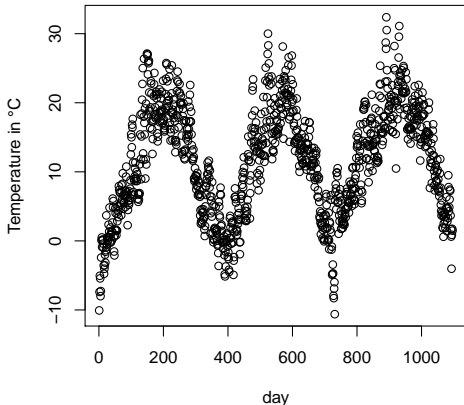
Example Time Series

- Temperature in Hamburg every day at 12:00
- Three years of data (1980, 1996, 2014)

```

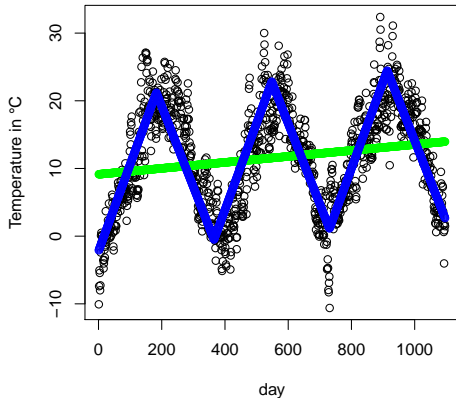
1 d = read.csv("temp-hamburg.csv", header=TRUE)
2 d$Lat = NULL; d$Lon = NULL
3 colnames(d) = c("h", "t")
4 d$t = d$t - 273.15 # convert degree Kelvin to
   ↪ Celcius
5 plot(d$t, xlab="day", ylab="Temperature in C")
6
7 pdf("hamburg-temp-models.pdf", width=5,height=5)
8 plot(d$t, xlab="day", ylab="Temperature in C")
9
10 # Enumerate values
11 d$index=1:nrow(d)
12
13 # General trend
14 m = lm("t ~ index", data=d)
15 points(predict(m, d), col="green")
16
17 # Summer/Winter model per day of the year
18 d$day=c(rep(c(1:183, 182:1),3),0)
19 m = lm("t ~ day + index", data=d)
20 points(predict(m, d), col="blue");
21
22 library(forecast)
23 # Convert data to a time series
24 ts = ts(d$t, frequency = 365, start= c(1980, 1))
25 # Apply a model for non-seasonal data on
   ↪ seasonal adjusted data
26 tsmod = stlm(ts, modelfunction=ar)
27 plot(forecast(tsmod, h=365))

```



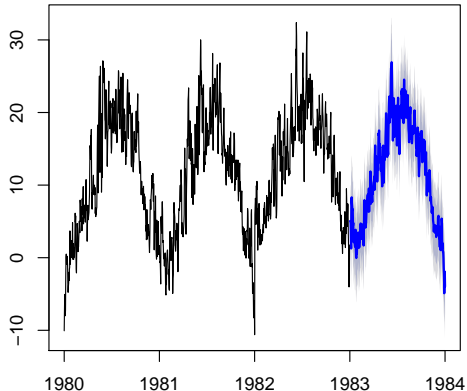
Temperature time series

Example Time Series Models



Linear models for trend & winter/summer cycle

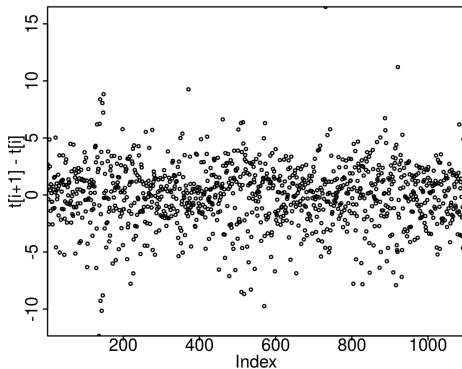
Forecasts from STL + AR(14)



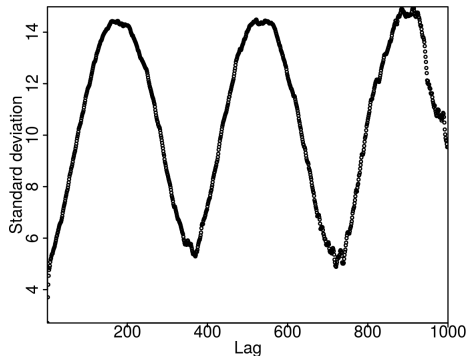
Using the forecast package/stlm()

Lag

- Lag: fixed time displacement, e.g., X_5 and X_2 have lag 3
- Lag operator returns the previous element ($LX_t = X_{t-1}$) [4]
- Plots help to identify patterns and determine if timeseries is random



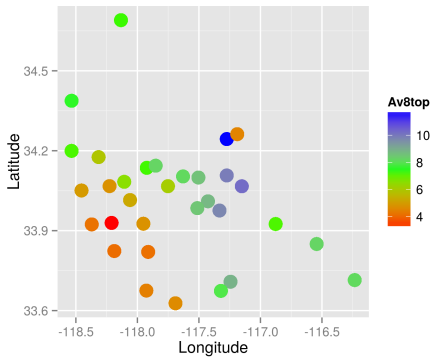
Lag plot, difference for Lag 1 for each element



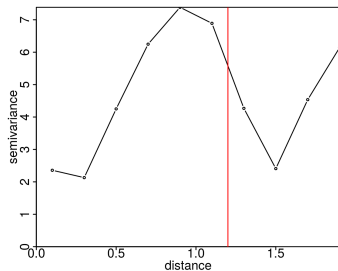
Standard deviation for all lags

Modeling Spatial Data

- Assume we measure data (z) on a surface (x,y), e.g., oil reservoirs
- **Kriging**: “methods are statistical estimation algorithms that curve-fit known point data to produce a predictive surface” [2, p.116]
- Variogram: variance between (x,y) pairs depending on the distance



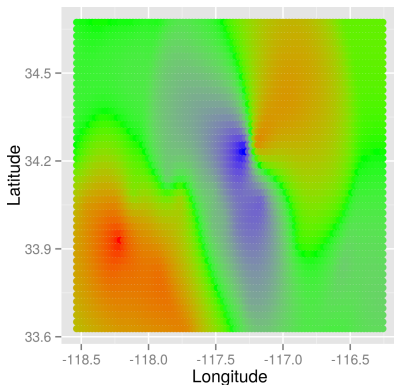
Stations and their observed ozone values



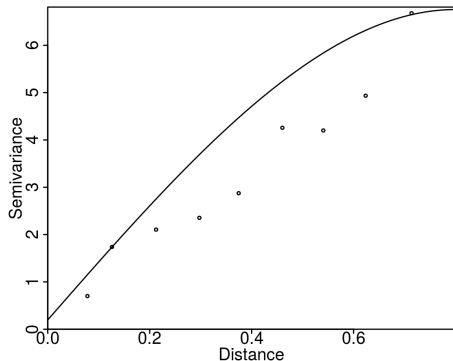
Estimated variogram with bins. A distance $>$ red line, only green stations are compared

Kriging with Different Methods

- Models: linear, exponential, Gaussian, spherical



Spherical model



Spherical model prediction and actual variance

Summary

- Statistics allows us to describe properties and infer properties
 - Sample: subset of data
 - Population: complete set of items that is subject to analysis
 - Independent vs. dependent variables
- Descriptive statistics helps analyzing samples
 - Univariate: 5-number summary, quantiles, histograms, boxplots, density
 - Bivariate: correlation of two variables
- Inductive statistics provide concepts for inferring knowledge
 - Aim: reject null-hypothesis
 - Careful investigation of model accuracy is needed
 - Methods have certain requirements to data distribution
 - Regression method with linear models
 - Time series: variables depend on previous state
 - Geospatial data is treated with Kriging
- More about prediction in the machine learning lecture

Bibliography

- 1 Book: Statistik macchiato – Cartoon-Stochastikkurs für Schüler und Studenten. Andreas Lindenberg, Irmgard Wagner. Pearson.
- 2 Book: Data Science for Dummies. Lillian Pierson. A Wiley Brand.
- 4 Lag operator, Wikipedia. https://en.wikipedia.org/wiki/Lag_operator
- 10 Wikipedia
- 21 Book: Y. Dodge. The Oxford Dictionary of Statistical Terms. ISBN 0-19-920613-9
- 22 https://en.wikipedia.org/wiki/Pearson_product-moment_correlation_coefficient
- 23 <http://blog.yhathq.com/posts/r-lm-summary.html>
- 24 Forecasting: principles and practise <https://www.otexts.org/fpp>