

```

1: ///////////////////////////////////////////////////////////////////
2: Baseline Routines
3:
4: maXbox Starter 84 - Baseline in Code - Max Kleiner
5: Machine Learning XII
6:
7: "In 1950, Alan Turing published a seminal paper titled "Computing
  Machinery and Intelligence" in Mind magazine. In this detailed paper the
  question "Can Machines Think?" was proposed. The paper suggested
  abandoning the quest to define if a machine can think, to instead test the
  machine with the 'imitation game'."
8:   - https://www.unite.ai/what-is-the-turing-test-and-why-does-it-matter/
9:
10: IEEE (IEEE Std. No. 610.12-1990) defines baseline as an agreed description
  and review of attributes of product, that afterward serve as basis for
  further development and defining change, and this changing can be done
  only through formal change control procedures".
11:
12: Software baseline instability has done more to undermine acquisition
  credibility and complicate effective management of the acquisition of
  software-intensive systems than the inability to establish realistic
  software development cost.
13: Our first baseline is synthetic data generated by a synthetic data random
  generator (sd-rage):
14:
15:   {
16:     "date": "2021-3-4",
17:     "confirmed": 36223,
18:     "deaths": 1483,
19:     "recovered": 33632,
20:     "location": EU
21:   }
22:
23: In order to be useful for a machine learning classifier, the synthetic
  data should have certain properties. While the data can be categorical,
  binary, or numerical, the length of the dataset should be arbitrary and
  the data should be randomly generated. The random processes used to
  generate the data should be controllable and based on various statistical
  distributions. Random noise may also be placed in the dataset.
24:
25: • Random Number Generator •
26:
27: A random number generator is a routine that produces a sequence of numbers
  that would pass statistical or probabilistic tests for randomness. To be
  really strict, routines that generate random numbers are said to be
  pseudorandom number generators (often abbreviated to PRNG) to
  differentiate them from true random number generators that rely on some
  kind of random events happening at a quantum level.
28:
29: type
30:   TStRandomBase = class
31:   private
32:   protected
33:     function rbMarsagliaGamma(aShape : double) : double;
34:     function rbMontyPythonNormal : double;
35:   public
36:     {uniform distributions}
37:     function AsFloat : double; virtual; abstract;
38:     function AsInt(aUpperLimit : integer) : integer;
39:     function AsIntInRange(aLowerLimit : integer;
40:                           aUpperLimit : integer) : integer;
41:

```

```

42:     {continuous non-uniform distributions}
43:     function AsBeta(aShapel, aShape2 : double) : double;
44:     function AsCauchy : double;
45:     function AsChiSquared(aFreedom : integer) : double;
46:     function AsErlang(aMean : double;
47:         aOrder : integer) : double;
48:     function AsExponential(aMean : double) : double;
49:     function AsF(aFreedom1 : integer;
50:         aFreedom2 : integer) : double;
51:     function AsGamma(aShape : double; aScale : double) : double;
52:     function AsLogNormal(aMean : double;
53:         aStdDev : double) : double;
54:     function AsNormal(aMean : double;
55:         aStdDev : double) : double;
56:     function AsT(aFreedom : integer) : double;
57:     function AsWeibull(aShape : double;
58:         aScale : double) : double;
59: end;

```

```

60:
61: TStRandomSystem = class(TStRandomBase)
62:     private
63:         FSeed : integer;
64:     protected
65:         procedure rsSetSeed(aValue : integer);
66:     public
67:         constructor Create(aSeed : integer);
68:         function AsFloat : double; override;
69:         property Seed : integer read FSeed write rsSetSeed;
70: end;

```

71:

72: <https://github.com/TurboPack/SysTools/blob/master/source/StRandom.pas>

73:

74: Often, **in** conversation, people use the term random when they really mean arbitrary. When one asks **for** an arbitrary number, one **is** saying that one doesn't really care what number one gets: almost any number will **do**. By contrast, a random number **is** a precisely defined mathematical concept: every number should be equally likely **to** occur. A random number will satisfy someone who needs an arbitrary number, but **not** the other way around.

75:

76: The first test **is** the simplest: the uniformity test. This **is** the one we were discussing earlier. Basically, the random numbers we generate are going **to** be checked **to** see that they uniformly cover the range **0.0 to 1.0**. We create **100** buckets, generate **1,000,000** random numbers, **and** slot them into each bucket. Bucket **0** gets all the random numbers from **0.0 to 0.01**, bucket **1** gets them from **0.01 to 0.02**, **and** so **on**. The probability **of** a random number falling into a particular bucket **is** obviously **0.01**. We calculate the chi-square value **for** our test **and** check that against the standard table, **using** the **99** degrees **of** freedom line.

77:

78:

```

79: SelfTestCRandom;
80: writeln('GetRandomSeedX: '+inttostr64(GetRandomSeedX));
81:
82: strand:= TStRandomSystem.create(42); //if 0 then randomseedx
83: writeln(floattostr(strand.asFloat))
84: writeln(floattostr(strand.asFloat))
85: writeln(floattostr({System.}Random(10)));
86: writeln(floattostr({System.}RandomE));
87: writeln(itoa(strand.asInt(10000)))
88: writeln(itoa(strand.asInt(10000)))
89: writeln(itoa(strand.asInt($E15)))

```

```
90:
91: UniformityTest2(strand, ChiSquare, DegrFreedom)
92: writeln('ChiSquare: '+floattostr(ChiSquare)+' ,
degsfreedom: '+itoa(degsfreedom));
93:
94: UniformityTest3(strand, ChiSquare, DegrFreedom)
95: writeln('ChiSquare: '+floattostr(ChiSquare)+' ,
degsfreedom: '+itoa(degsfreedom));
96:
97: strand.Free;
98:
99: Current theories state that quantum events are truly random. The time of
the decay of a radioactive atom into its byproducts cannot be predicted;
all we can say is that there is a certain probability that it will decay
within a certain period of time, and we can estimate that probability by
observing the decay of many, many atoms.
100:
101: • What is baseline model in machine learning? •
102:
103: A baseline is a method that uses heuristics, simple summary statistics,
randomness, or machine learning to create predictions for a dataset. You
can use these predictions to measure the baseline's performance (e.g.,
accuracy)- this metric will then become what you compare any other machine
learning algorithm against.
104:
105: The three most commonly used baseline algorithms are:
106:
107: •Random Prediction Algorithm.
108: •Zero Rule Algorithm.
109: •Naive Bayes Algorithm
110:
111: When starting on a new problem that is more sticky than a conventional
classification or regression problem, it is a good idea to first devise a
random prediction algorithm that is specific to your prediction problem.
Later you can improve upon this and devise a zero rule algorithm. Come up
with a baseline model. As you will see, Naive Bayes, acts as a
tremendously good baseline model because of its probabilistic approach. It
is comprehensible and simple.
112:
113: A machine learning algorithm tries to learn a function that models the
relationship between the input (feature) data and the target variable (or
label). When you test it, you will typically measure performance in one
way or another. For example, your algorithm maybe 85% accurate. But what
does this mean? You can infer this meaning by comparing it with a
baseline's performance.
114:
115: Typical baselines include those supported by scikit-learn's "dummy"
estimators:
116:
117: Classification baselines:
118:
119: • stratified": generates predictions by respecting the training set's
class distribution.
120: • most_frequent": always predicts the most frequent label in the
training set.
121: • prior": always predicts the class that maximizes the class prior.
122: • uniform": generates predictions uniformly at random.
123: • constant": always predicts a constant label that is provided by the
user.
124:
125: This is useful for metrics that evaluate a non-majority class.
126:
```

```

127: • Conclusion:
128: Hence, more often than not I build a Naive Bayes model as my baseline and
    then go on towards building more complex ones such as Decision Tree,
    Support Vector Machine or Random Forest.
129:
130: If you are dealing with a specific domain of machine learning (such as
    recommender systems), then you will typically pick baselines that are
    current state-of-the-art (SoTA) approaches - since you will usually want to
    demonstrate that your approach does better than these. For example, while
    you evaluate a new collaborative filtering algorithm, you may want to
    compare it to matrix factorization - which itself is a learning algorithm,
    but is now a popular baseline since it has been so successful in
    recommender system research.
131:
132:
133: Appendix for Python to get a Polynomial Regression Baseline:
134:
135: Base Schema for Polynomial Regression
136: -----
137: A straightforward way of doing multivariate polynomial regression for
    Python.
138:
139: import numpy as np
140: import matplotlib.pyplot as plt
141: from sklearn.preprocessing import PolynomialFeatures
142: from sklearn.linear_model import LinearRegression
143: from sklearn.pipeline import make_pipeline
144:
145: rng=np.random.RandomState(1)
146: # get some sin data
147: X=10 * rng.rand(50)
148: y=np.sin(X)+0.1*rng.randn(50)
149:
150: poly=make_pipeline(PolynomialFeatures(degree=7),LinearRegression())
151: Xfit=np.linspace(0,10,1000)
152: # train and predict
153: poly.fit(X[:,np.newaxis],y)
154: >>>Pipeline(steps=[('polynomialfeatures', PolynomialFeatures(degree=7)),
155:                    ('linearregression', LinearRegression())])
156: yfit=poly.predict(Xfit[:,np.newaxis])
157: # plot it
158: plt.scatter(X,y)
159: >>><matplotlib.collections.PathCollection object at 0x000000523600E470>
160: plt.title("Polynomial mX7fit graph")
161: >>>Text(0.5, 1.0, 'Polynomial mX7fit graph')
162: plt.plot(Xfit, yfit)
163: >>>[<matplotlib.lines.Line2D object at 0x0000005236C1CF28>]
164: plt.show()
165:
166: Ref:
167:     https://simulatoran.com/what-is-baseline-model-in-machine-learning/
168:     https://pomber.github.io/covid19/timeseries.json
169:     script: 1026_json_automation_refactor2.txt
170: Doc:
171:     https://maxbox4.wordpress.com
172: >>> https://entwickler-konferenz.de/speaker/max-kleiner/
173:

```