

Nearest Neighbor Correlation in production of Airplane-Induced Pollution

George Durendal

Abstract

We present and discuss a nearest neighbors correlation in the production of airplane-induced air pollution. A k-nearest neighbors (k-NN) regression consistently produced a higher r^2 values than multiple other types of regressions as a function of time. There are therefore clusters of airplane-induced air pollution high-production and low-production days. Private planes are a potential source of lead pollution, so this can have a major impact on public policy.

1 Research

An AutoML regression analysis of data closely correlated with the amount of airplane-induced air pollution produced by airlines consistently yielded the result that k-nearest neighbors (k-NN) regression consistently produced a higher r^2 values than multiple other types of regressions. The regressions were plotted as a function of day index in the year, so Jan 1st is 1, Dec 31st is either 365 or 366 during a leap year. The data is based on U.S. composite federal data, not localized to any particular region of the US.

This suggests that high airplane-induced air pollution (ApIA pollution) on one particular day is correlated with high levels of ApIA pollution produced on nearby days. We note that this is produced ApIA pollution, i.e. it makes no statement on lingering pollution. Hence the nearest day correlation is not due to lingering pollution from nearby dates.

We posit that this nearest neighbor correlation is due to popular travel dates such as holidays, peak business dates, weekends, and peak seasonal travel dates unrelated to common holidays (e.g. summer). Weekends introduce a common and consistent occurrence of

One competing algorithm is NuSVR, which beat k-NN in some tests, although never reaching as

high a r^2 as those achieved by some k-NN regressions. NuSVR is here defined as Support Vector Regression (SVR) with a parameter Nu controlling the number of support vectors.

2 Discussion

One notable point is that the k-NN beats regression algorithms of numerous types. This implies that, in regard to production of ApIA pollution, production levels from nearby dates are a significant indicator of current production levels. Moreover, this nearest neighbor relationship is strong enough to outperform other regressions attempted. For instance, historically the US sees elevated travel levels at the end of the year, during its major holiday seasons. This positive trend was not strong enough to make some regressor which might show some positive trend towards the tail end of the dataset outperform the k-NN regressor. Thus, in terms of air travel and associated ApIA pollution, the nearest neighbor relation seems to be more important than any A natural corollary of this nearest neighbor relation is that ApIA pollution levels are concentrated around specific clusters of dates, particularly holidays and high-travel dates. Lingering ApIA pollution is only compounded by elevated production levels on adjacent dates. The data used here is generalized to that of the entire US, so it makes no prediction about ApIA levels in any specific region or locality. Although we may naturally posit that ApIA pollution levels are highest around public and private airports and the areas where those planes fly over or where the pollution from them disperses to. An additional corollary is that ApIA production is lowest during clusters of dates with lower levels of airplane travel. Given the fact that private aircraft are known to emit toxins such as lead or lead composites, this may be applicable to public policy. It might, for instance, be prudent to sched-

3 Regression Plots

Model	R-Squared (approx.)
KNeighborsRegressor	0.40
SVR	0.35
LinearSVR	0.32
LinearRegression	0.30
TransformedTargetRegressor	0.28
LassoLarsCV	0.27
LassoLars	0.26
OrthogonalMatchingPursuit	0.25
Lasso	0.24
Lars	0.23
Ridge	0.22
LassoCV	0.21
ElasticNet	0.20
HuberRegressor	0.19
SVR	0.18
BayesianRidge	0.17
RandomForestRegressor	0.16
RandomForestRegressor	0.15
RidgeCV	0.14
ExtraTreesRegressor	0.13
ExtraTreesRegressor	0.12
XGBRegressor	0.11
GradientBoostingRegressor	0.10
DecisionTreeRegressor	0.09
DecisionTreeRegressor	0.08
RANSACRegressor	0.07
LassoLars	0.06
ElasticNet	0.05
DummyRegressor	0.04
HasGradientBoostingRegressor	0.03
GBMRegressor	0.02
GaussianM	0.01
MLPRegressor	0.01
SGDRegressor	0.01
PassiveAggressiveRegressor	0.01

Model	R-Squared (approx.)
KNeighborsRegressor	0.38
GaussianNB	0.36
NaiveBayes	0.34
RandomForestClassifier	0.32
AdaBoostClassifier	0.30
ExtraTreeClassifier	0.28
LassoCV	0.26
ElasticNet	0.24
Bayesian Ridge	0.22
LogisticRegression	0.20
LinearRegression	0.18
LarsCV	0.16
OrthogonalMatchingPursuit	0.14
LassoLarsCV	0.12
TransformedTargetRegressor	0.10
Lasso, SVR	0.08
BaggingRegressor	0.06
RandomizedGradientRegressor	0.04
DecisionTreeRegressor	0.02
ExtraTreeRegressor	0.01
GradientBoostingRegressor	0.00
DecisionTreeRegressor	0.00
LogisticRegression	0.00
ExtraTreeRegressor	0.00
LogisticRegression	0.00
HistGradientBoostingRegressor	0.00
Lasso	0.00
ElasticNet	0.00
DummyRegressor	0.00
SVC	0.00
SGDRegressor	0.00
PassiveAggressiveRegressor	0.00
GaussianProcessRegressor	0.00

Model	R-Squared (approx.)
KNeighborsRegressor	0.45
AdaBoostRegressor	0.42
SVR	0.40
AdaBoostRegressor	0.35
Ridge	0.32
RandomForestRegressor	0.30
BayesianRidge	0.28
ElasticNetCV	0.25
LassoCV	0.24
LinearRegression	0.23
TransformedTargetRegressor	0.22
Lars	0.21
LarsCV	0.20
OrthogonalMatchingPursuit	0.19
LassoLarsCV	0.18
LassoLarsIC	0.17
HuberRegressor	0.16
LinearSVR	0.15
RandomForestRegressor	0.14
BaggingRegressor	0.13
ExtraTreesRegressor	0.12
RANSACRegressor	0.11
ExtraTreeRegressor	0.10
GradientBoostingRegressor	0.09
DecisionTreeRegressor	0.08
XGBRegressor	0.07
HistGradientBoostingRegressor	0.06
LassoLars	0.05
ElasticNet	0.04
DummyRegressor	0.03
LogisticRegression	0.02
MLPRegressor	0.01
PassiveAggressiveRegressor	0.00
GaussianProcessRegressor	0.00

Model	R-Squared
SVR	~0.48
KNeighborsRegressor	~0.46
ElasticNetCV	~0.44
LassoCV	~0.42
BayesianRidge	~0.40
RandomForestClassifier	~0.39
AdaBoostClassifier	~0.37
LinearRegression	~0.35
LancVCV	~0.33
OtpogonalMatchingPursuit	~0.31
LassoLncVCV	~0.29
TransformedTargetRegressor	~0.27
HuberRegressor	~0.25
NUSVR	~0.23
RANSACRegressor	~0.21
BaggingRegressor	~0.19
RandomForestRegressor	~0.17
LinearSVR	~0.15
AdaBoostRegressor	~0.13
ExtraTreeRegressor	~0.11
PassiveAgressiveRegressor	~0.09
XGBRegressor	~0.07
GradientBoostingRegressor	~0.05
LassoLars	~0.03
Lasso	~0.02
ElasticNet	~0.01
DummyRegressor	~0.00
HuberRegressor	~0.00
HasGradientBoostingRegressor	~0.00
DecisionTreeRegressor	~0.00
GridSearchCV	~0.00
SQRRegressor	~0.00
MLPRegressor	~0.00
GaussianProcessRegressor	~0.00
KernalRidge	~0.00

Model	R Squared (approx.)
K-Nearest Regressor	0.48
Linear	0.38
Lasso	0.35
LassoLars	0.35
Orthogonal Matching Pursuit	0.35
Transformative Boosting	0.35
Linear Regression	0.35
NuSVR	0.35
Ridge	0.35
Huber Ridge	0.35
BayesAimRidge	0.30
SVR	0.28
RidgeCV	0.22
Bagging Regressor	0.20
Random Forest Regressor	0.18
AdaBoost Regressor	0.18
Extra Tree Regressor	0.08
LGBM Regressor	0.05
XGBoost Regressor	0.05
Hist Gradient Boosting Regressor	0.05
LassoCV	0.05
Dummy Regressor	0.05
Elastic Net	0.05
Lasso	0.05
LassoLars	0.05
LassoCV	0.05
Gradient Boosting Regressor	0.05
XGB Regressor	0.05
Decision Tree Regressor	0.05
MLP Regressor	0.05
SGD Regressor	0.05
Gaussian Process Regressor	0.05
Passive Aggressive Regressor	0.05
Kernel Ridge	0.05

Model	R-Squared (approx.)
NuSVR	0.45
KNeighborsClassifier	0.42
OrthogonalMatchingPursuit	0.40
LassoLarsCV	0.38
LinearRegression	0.37
LassoLars	0.36
LassoLarsIC	0.35
TransformedTargetRegressor	0.34
Ridge	0.33
ElasticNetCV	0.32
BayesianRidge	0.31
LinearSVR	0.30
BaggingRegressor	0.29
RandomForestRegressor	0.28
ExtraTreesRegressor	0.27
PassiveAggressiveRegressor	0.26
GradientBoostingRegressor	0.25
XGBRegressor	0.24
DecisionTreeRegressor	0.23
RANSACRegressor	0.22
LGBMRegressor	0.21
LassoLars	0.20
ElasticNet	0.19
DummyRegressor	0.18
HistGradientBoostingRegressor	0.17
AdaBoostRegressor	0.16
SGDRegressor	0.15
MLPRegressor	0.14
GaussianProcessRegressor	0.13
KNeighborsRegressor	0.12

Model	R-Squared (approx.)
KNeighborhoodRegressor	0.45
PassiveAggressiveRegressor	0.38
LassoLarsCV	0.35
LarsCV	0.32
Lars	0.30
LassoLarsIC	0.28
OrthogonalMatchingPursuit	0.28
LinearRegression	0.28
TransformedTargetRegressor	0.28
Ridge	0.28
RANSACRegressor	0.28
Bayesianridge	0.28
HuberRegressor	0.28
LassoCV	0.28
ElasticNetCV	0.28
LinearSVR	0.28
RandomForestRegressor	0.25
NaiveBayes	0.22
ExtraTreeRegressor	0.18
AdaBoostRegressor	0.15
GradientBoostingRegressor	0.12
LinearSVC	0.10
Lasso	0.08
ElasticNet	0.05
DummyRegressor	0.05
GradientBoostingRegressor	0.05
DecisionTreeRegressor	0.05
AdaBoostClassifier	0.05
MLPRegressor	0.05
ExtraTreeRegressor	0.05
SGDRegressor	0.05
GaussianProcessRegressor	0.05
KNeighborsRegressor	0.05

Figure 8: Plot 8

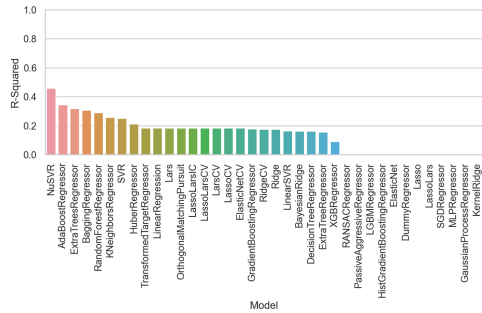


Figure 9: Plot 9

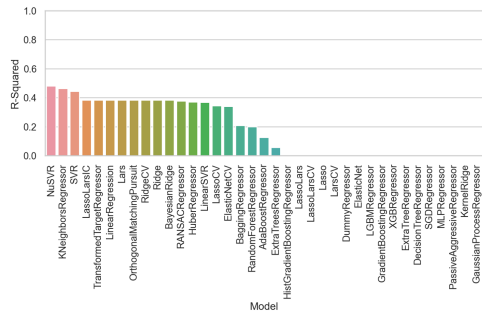


Figure 10: Plot 10

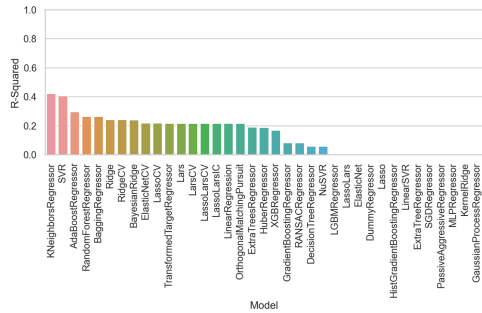


Figure 11: Plot 11

