



WHY GOOD MODELS FAIL

HOW AMBIGUITY, UNCERTAINTY, AND BAD SCIENCE
CAUSE MOST DATA SCIENCE PROJECTS TO FAIL





87%

OF ALL DATA SCIENCE PROJECTS FAIL
-VENTUREBEAT



The background features decorative 3D geometric shapes in the corners. In the top-left and top-right corners, there are stacks of light green rectangular blocks. In the bottom-right corner, there is a solid dark green cube with a shadow. In the bottom-center, there are light green L-shaped geometric outlines.

REASON ONE: **AMBIGUITY**

WHEN REQUIREMENTS ARE HARD TO DEFINE



AMBIGUITY OCCURS WHEN WE DON'T KNOW WHAT QUESTIONS TO ASK





HOW TO NAVIGATE **AMBIGUITY** IN DATA SCIENCE?





1. FILL KNOWLEDGE GAPS

BECOME A SUBJECT MATTER EXPERT



The top-left and top-right corners of the slide feature decorative elements made of light green 3D rectangular blocks. The top-left corner has a stack of five blocks of varying sizes, while the top-right corner has a smaller stack of three blocks.

2. ITERATE QUICKLY

RESEARCH, PRESENT, GET FEEDBACK





3. DEFER TO THE EXPERTS

LET THEM DRIVE THE RESEARCH





4. EXPLORE WITHOUT BUILDING

UNDERSTAND THE PROBLEM FULLY, THEN START DESIGNING



The background features isometric geometric shapes in the corners. Top-left: A stack of four rectangular blocks. Top-right: A stack of three rectangular blocks. Bottom-left: A large L-shaped block. Bottom-right: A solid green cube with a shadow.

REASON TWO: UNCERTAINTY

WHEN SUCCESS IS HARD TO DEFINE



HOW TO NAVIGATE UNCERTAINTY IN DATA SCIENCE?



1. UNDERSTAND THE PROCESS

THE FULL LIFECYCLE OF A RECORD





2. UNDERSTAND THE KPIS

THE SOURCE, TARGET, CALCULATION, STAKEHOLDERS





3. UNDERSTAND THE DOLLARS

SAVINGS/REVENUE CALCULATIONS MUST BE DONE METHODICALLY





**THE PROCESS IMPACTS THE KPI'S
WHICH IMPACT THE DOLLARS**



REASON THREE: BAD SCIENCE

WHEN VALUE IS MIS-MEASURED





1. INCORRECT METRICS

INAPPROPRIATE METRICS ARE THE HALLMARK OF BAD SCIENCE





2. BAD STATISTICS

THE MOST UNDER-APPRECIATED ASPECT OF DATA SCIENCE





BIGGEST STATISTICS MISTAKES IN DATA SCIENCE



Jessie Lamontagne · 1st

Machine Learning | Data Science | Analytics

2mo (edited) ...

Describing only the most vanilla undergrad-level regression as Regression, looking to graph neural networks as the only way to get causal inference while ignoring decades of research on causal inference with tabular data, treating point estimates as truth, not calculating uncertainty at all "because big data".

Celebrate ·



23

Reply ·

10 Replies

NOT STATISTICS MISTAKES IN DATA SCIENCE





Jessie Lamontagne · 1st

Machine Learning | Data Science | Analytics

Describing only the most vanilla undergrad-level regression as Regression, looking to graph neural networks as the only way to get causal inference while ignoring decades of research on causal inference with tabular data, treating point estimates as truth, not calculating uncertainty at all "because big data".

Celebrate ·



23

Reply ·

10 Replies

2mo (edited) ...



Leigh McCormack · 1st

Proven Builder of Innovative Teams, Products, and Processes | Hea...

Leaning too heavily on parametric tests when most real life data don't fit those molds!

Insightful ·



13

Reply ·

2 Replies

2mo ...

10 STATISTICS MISTAKES IN DATA SCIENCE





Jessie Lamontagne · 1st

Machine Learning | Data Science | Analytics

Describing only the most vanilla undergrad-level regression as Regression, looking to graph neural networks as the only way to get causal inference while ignoring decades of research on causal inference with tabular data, treating point estimates as truth, not calculating uncertainty at all "because big data".

Celebrate ·



23

Reply ·

10 Replies

2mo (edited) ...



Leigh McCormack · 1st

Proven Builder of Innovative Teams, Products, and Processes | Hea...

Leaning too heavily on parametric tests when most real life data don't fit those molds!

Insightful ·



13

Reply ·

2 Replies

2mo ...

NOT STATISTICS MISTAKES IN DATA SCIENCE

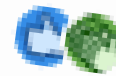


Rodrigo Rivera · 1st

🔧 Data Products for Business Users Builder | 🐱 Enabling Anyone t...

Basing all decisions around the p-value

Like ·



9

Reply ·

3 Replies

2mo ...



Jessie Lamontagne · 1st

Machine Learning | Data Science | Analytics

Describing only the most vanilla undergrad-level regression as Regression, looking to graph neural networks as the only way to get causal inference while ignoring decades of research on causal inference with tabular data, treating point estimates as truth, not calculating uncertainty at all "because big data".

Celebrate · 🥳🌱 23 | Reply · 10 Replies

2mo (edited) ...



Leigh McCormack · 1st

Proven Builder of Innovative Teams, Products, and Processes | Hea...

Leaning too heavily on parametric tests when most real life data don't fit those molds!

Insightful · 🗨️🌱 13 | Reply · 2 Replies

2mo ...

NOT STATISTICS MISTAKES IN DATA SCIENCE



Kristen Kehrher (She/Her) · 1st

Data Scientist | Developer Advocate @ CometML | The Cool Data P...

Not so much fellow statisticians, but I was a broken record saying "correlation is not causation we need to test these results" to stakeholders after sharing model output.

Like · 🗨️🌱 9 | Reply · 1 Reply

2mo ...



Rodrigo Rivera · 1st

🔧 Data Products for Business Users Builder | 🐱 Enabling Anyone t...

Basing all decisions around the p-value

Like · 🗨️🌱 9 | Reply · 3 Replies

2mo ...



Jessie Lamontagne · 1st

Machine Learning | Data Science | Analytics

Describing only the most vanilla undergrad-level regression as Regression, looking to graph neural networks as the only way to get causal inference while ignoring decades of research on causal inference with tabular data, treating point estimates as truth, not calculating uncertainty at all "because big data".

Celebrate



23

Reply

10 Replies

2mo (edited) ...



Leigh McCormack · 1st

Proven Builder of Innovative Teams, Products, and Processes | Hea...

Leaning too heavily on parametric tests when most real life data don't fit those molds!

Insightful

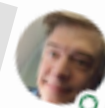


13

Reply

2 Replies

2mo ...



Peter Mancini (He/Him) · 1st

Principal at Stealth Startup

1. Not knowing what distribution their data is and using tools that depend upon a specific distribution, usually the normal distribution.
2. Relying on normal distribution tools de facto. Better to rely on tools that will work with the distribution of the data that you have.
3. No concept of uncertainty. I think this is a failure that goes through a much wider swathe of practitioners including fairly seasoned practitioners. In fact, all of the faults that you mentioned contain uncertainty as a key feature. You were going to be uncertain about every single outcome, sometimes a lot, sometimes a miniscule amount, usually something in between. Knowing it will help.

Celebrate



3

Reply

1 Reply

2mo ...



Kristen Kehrher (She/Her) · 1st

Data Scientist | Developer Advocate @ CometML | The Cool Data P...

Not so much fellow statisticians, but I was a broken record saying "correlation is not causation we need to test these results" to stakeholders after sharing model output.

Like



9

Reply

1 Reply

2mo ...

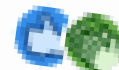


Rodrigo Rivera · 1st

Data Products for Business Users Builder | Enabling Anyone t...

Basing all decisions around the p-value

Like



9

Reply

3 Replies

mo ...

NOT STATISTICS M IN DATA SCIENCE

NOT STATISTICS M IN DATA SCIENCE

Jessie Lamontagne · 1st
Machine Learning | Data Science | Analytics
Describing only the most vanilla undergrad-level regression as Regression, looking to graph neural networks as the only way to get causal inference while ignoring decades of research on causal inference with tabular data, treating point estimates as truth, not calculating uncertainty at all "because big data".

Celebrate · 🥳🌍 23 | Reply · 10 Replies
2mo (edited) ...

Leigh McCormack · 1st
Proven Builder of Innovative Teams, Products, and Processes | Hea...
Leaning too heavily on parametric tests when most real life data don't fit those molds!

Insightful · 🗨️ 13 | Reply · 2 Replies
2mo ...

Peter Mancini (He/Him) · 1st
Principal at Stealth Startup
1. Not knowing what distribution their data is and using tools that depend upon a specific distribution, usually the normal distribution.
2. Relying on normal distribution tools de facto. Better to rely on tools that will work with the distribution of the data that you have.
3. No concept of uncertainty. I think this is a failure that goes through a much wider swathe of practitioners including fairly seasoned practitioners. In fact, all of the faults that you mentioned contain uncertainty as a key feature. You were going to be uncertain about every single outcome, sometimes a lot, sometimes a miniscule amount, usually something in between. Knowing it will help.

Celebrate · 🥳🌍 3 | Reply · 1 Reply
2mo ...

Kristen Kehrner (She/Her) · 1st
Data Scientist | Developer Advocate @ CometML | The Cool Data P...
Not so much fellow stat "correlation is not causation" but I was a broken record saying stakeholders after sharing results" to

Like · 🗨️ 9 | Reply
2mo ...


Joe Wyer (He/Him) · 1st
Head of Economics and Applied Science @ Haus | We're Hiring!
Thinking statistical significance should be the filter for go/no go of business decisions.

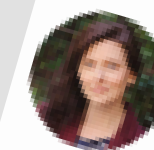
Like · 🗨️ 4 | Reply · 2 Replies
1mo ...


Rodrigo Rivera · 1st
Data Products for Business Users Builder | 🐱 Enabling Anyone t...
around the p-value

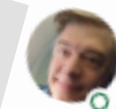
Reply · 3 Replies
mo ...


NOT STATISTICS M IN DATA SCIENCE


 **Zion Pibowei** · 1st
Principal Data Scientist | Head of Data Science @Periculum | Data ...
Using SMOTE when they should not 🤔
Celebrate · 🌍 2 | Reply · 1 Reply


 **Leigh McCormack** · 1st
Proven Builder of Innovative Teams, Products, and Processes | Hea...
Leaning too heavily on parametric tests when most real life data
don't fit those molds!
Insightful · 🗣️ 13 | Reply · 2 Replies

 **Jessie Lamontagne** · 1st
Machine Learning | Data Science | Analytics
Describing only the most vanilla undergrad-level regression
Regression, looking to graph neural networks as the only way to
causal inference while ignoring decades of research on causal
inference with tabular data, treating point estimates as truth, not
calculating uncertainty at all "because big data".
Celebrate · 🗣️ 23 | Reply · 10 Replies

 **Peter Mancini (He/Him)** · 1st
Principal at Stealth Startup
1. Not knowing what distribution their data is and using tools that
depend upon a specific distribution, usually the normal distribution.
2. Relying on normal distribution tools de facto. Better to rely on
tools that will work with the distribution of the data that you have.
3. No concept of uncertainty. I think this is a failure that goes
through a much wider swathe of practitioners including fairly
seasoned practitioners. In fact, all of the faults that you mentioned
contain uncertainty as a key feature. You were going to be
uncertain about every single outcome, sometimes a lot, sometimes
a miniscule amount, usually something in between. Knowing it will
help.
Celebrate · 🗣️ 3 | Reply · 1 Reply

 **Kristen Kehrher (She/Her)** · 1st
Data Scientist | Developer Advocate @ CometML | The Cool Data P...
Not so much fellow stat
"correlation is not caus
stakeholders after shar
but I was a broken record saying
results" to
Like · 🗣️ 9 | Reply

 **Rodrigo Rivera** · 1st
Data Products for Business Users Builder | 🐱 Enabling Anyone t...
around the p-value
y · 3 Replies

 **Joe Wyer (He/Him)** · 1st
Head of Economics and Applied Science @ Haus | We're Hiring!
Thinking statistical significance should be the filter for go/no go of
business decisions.
Like · 🗣️ 4 | Reply · 2 Replies



Zion Pibowei · 1st

Principal Data Scientist | Head of Data Science @Periculum | Data ...

Using SMOTE when they should not 🤔

2mo ...

Celebrate · 🌍 2 | Reply · 1 Reply



Leigh McCormack · 1st

Proven Builder of Innovative Teams, Products, and Processes | Hea...

Leaning too heavily on parametric tests when most real life data don't fit those molds!

2mo ...

Insightful · 🌍 13 | Reply · 2 Replies



Jessie Lamontagne · 1st

Machine Learning | Data Science | Analytics

Describing only the most vanilla undergrad-level regression, looking to graph neural networks as the only way to causal inference while ignoring decades of research on causal inference with tabular data, treating point estimates as truth, not calculating uncertainty at all "because big data".

Celebrate · 🌍 23 | Reply · 10 Replies



Hai Tran (He/Him) · 1st

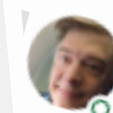
Data Science and Analytics Specialist

Successfully rejecting or failing to reject Null Hypothesis H_0 does not prove that the Alternative Hypothesis is True or False. We can only say whether H_1 is supported by the observed data or not.

Another common pitfall should be being mistaken for the correlation between 2 variables and thinking that it is a causation relationship between these.

2mo ...

Celebrate · 🌍 4 | Reply · 1 Reply



Peter Mancini (He/Him) · 1st

Principal at Stealth Startup

1. Not knowing what distribution their data is and using tools that depend upon a specific distribution, usually the normal distribution.
2. Relying on normal distribution tools de facto. Better to rely on tools that will work with the distribution of the data that you have.
3. No concept of uncertainty. I think this is a failure that goes through a much wider swathe of practitioners including fairly seasoned practitioners. In fact, all of the faults that you mentioned contain uncertainty as a key feature. You were going to be uncertain about every single outcome, sometimes a lot, sometimes a miniscule amount, usually something in between. Knowing it will help.

2mo ...

Celebrate · 🌍 3 | Reply · 1 Reply



Kristen Kehrher (She/Her) · 1st

Data Scientist | Developer Advocate @ CometML | The Cool Data P...

Not so much fellow stat "correlation is not causation" but I was a broken record saying stakeholders after sharing results" to

2mo ...

Like · 🌍 9 | Reply



Joe Wyer (He/Him) · 1st

Head of Economics and Applied Science @ Haus | We're Hiring!

Thinking statistical significance should be the filter for go/no go of business decisions.

1mo ...

Like · 🌍 4 | Reply · 2 Replies

Rodrigo Rivera · 1st

Data Products for Business Users Builder | 🐱 Enabling Anyone t...

around the p-value

mo ...

· 3 Replies



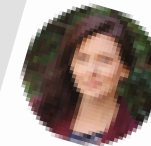
Zion Pibowei · 1st

Principal Data Scientist | Head of Data Science @Periculum | Data ...

Using SMOTE when they should not 🤔

2mo ...

Celebrate · 🌍 2 | Reply · 1 Reply



Leigh McCormack · 1st

Proven Builder of Innovative Teams, Products, and Processes | Hea...
Leaning too heavily on parametric tests when most real life data don't fit those molds!

2mo ...

Insightful · 🌍 13 | Reply · 2 Replies



Jessie Lamontagne · 1st

Machine Learning | Data Science | Analytics
Describing only the most vanilla undergrad-level regression models as the only way to understand the world is looking to graph neural networks as the only way to understand the world. Decades of research on causal inference as truth, not just correlation.

2mo (edited) ...



Aaron Sheldon · Following
VP of Engineering at Genstate

In general across all scientific fields it is failing to understand the strong limitations of the mathematical assumptions of statistical models. Basically applying statistical testing as a form of exploratory analysis. Or more bluntly assuming regression and its generalizations all the way up to deep learning are a magical oracle that can tell you how your system works.

Statistical tests and models assume a priori that you have a theory of how the system your are studying works. They simply estimate the quantitative bounds of your theory. Not generally whether it is right or wrong. This is even true of classical hypothesis testing.

Celebrate · 🌍 34 | Reply · 4 Replies



Kristen Kehrer (She/Her) · 1st

Data Scientist | Developer Advocate @ CometML | The Cool Data P...
Not so much fellow stat but I was a broken record saying "correlation is not causation" to stakeholders after sharing results to

2mo ...

Like · 🌍 9 | Reply



Joe Wyer (He/Him) · 1st

Head of Economics and Applied Science @ Haus | We're Hiring!

Thinking statistical significance should be the filter for go/no go of business decisions.

1mo ...

Like · 🌍 4 | Reply · 2 Replies

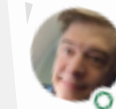


Rodrigo Rivera · 1st

Data Products for Business Users Builder | 🐱 Enabling Anyone t...
Racism around the p-value

mo ...

· 3 Replies



Peter Mancini (He/Him) · 1st

Principal at Stealth Startup

1. Not knowing what distribution their data is and using tools that depend upon a specific distribution, usually the normal distribution.
2. Relying on normal distribution tools de facto. Better to rely on tools that will work with the distribution of the data that you have.

3. No concept of uncertainty. I think this is a failure that goes through a much wider swathe of practitioners including fairly seasoned practitioners. In fact, all of the faults that you mentioned contain uncertainty as a key feature. You were going to be uncertain about every single outcome, sometimes a lot, sometimes a miniscule amount, usually something in between. Knowing it will help.

2mo ...

Celebrate · 🌍 3 | Reply · 1 Reply



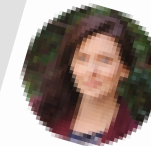
Zion Pibowei · 1st

Principal Data Scientist | Head of Data Science @Periculum | Data ...

Using SMOTE when they should not 🤔

2mo ...

Celebrate · 🌍 2 | Reply · 1 Reply



Leigh McCormack · 1st

Proven Builder of Innovative Teams, Products, and Processes | Hea...
Leaning too heavily on parametric tests when most real life data don't fit those molds!

2mo ...

Insightful · 🗣️ 13 | Reply · 2 Replies



Jessie Lamontagne · 1st

Machine Learning | Data Science | Analytics
Describing only the most vanilla undergrad-level regression models as the only way to understand the world is looking to graph neural networks as the only way to understand the world. Decades of research on causal inference as truth, not just correlation.

2mo (edited) ...



Aaron Sheldon · Following
VP of Engineering at Genstate

In general across all scientific fields it is failing to understand the strong limitations of the mathematical assumptions of statistical models. Basically applying statistical testing as a form of exploratory analysis. Or more bluntly assuming regression and its generalizations all the way up to deep learning are a magical oracle that can tell you how your system works.

Statistical tests and models assume a priori that you have a theory of how the system your are studying works. They simply estimate the quantitative bounds of your theory. Not generally whether it is right or wrong. This is even true of classical hypothesis testing.

Celebrate · 🌍 34 | Reply · 4 Replies



Kristen Kehr (She/Her) · 1st

Data Scientist | Developer Advocate @ CometML | The Cool Data P...
... but I was a broken record saying "correlation is not causation" to stakeholders after sharing results to

2mo ...

Like · 🗣️ 9 | Reply

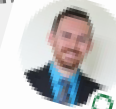


Hai Tran (He/Him) · 1st
Data Science and Analytics Specialist

Successfully rejecting or failing to reject Null Hypothesis H_0 does not prove that the Alternative Hypothesis is True or False. We can only say whether H_1 is supported by the observed data or not.

Another common pitfall should be being mistaken for the correlation between 2 variables and thinking that it is a causal relationship between these.

Celebrate · 🌍 4 | Reply · 1 Reply



Tyler Buffington, PhD · 1st
Senior Data Scientist

1. Not understanding decision analysis and value of information (VOI). For example, building an extremely complicated model when a simple analysis is sufficient to guide the business to the correct decision.
2. Failing to understand that underpowered experiments are untrustworthy even if the result is statistically significant.

2mo ...

Insightful · 🗣️ 6 | Reply · 2 Replies



Rodrigo Rivera · 1st

Data Products for Business Users Builder | 🐱 Enabling Anyone t...
... around the p-value

... 3 Replies



Joe Wyer (He/Him) · 1st

Head of Economics and Applied Science @ Haus | We're Hiring!

Thinking statistical significance should be the filter for go/no go of business decisions.

1mo ...

Like · 🗣️ 4 | Reply · 2 Replies



Zion Pibowei · 1st

Principal Data Scientist | Head of Data Science @Periculum | Data ...

Using SMOTE when they should not 🤔

2mo ...

Celebrate · 🌍 2

Reply · 1

2mo ...



Leigh McCormack · 1st

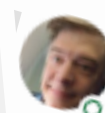
Proven Builder of Innovative Teams, Products, and Processes | Hea...
Leaning too heavily on parametric tests when most real life data don't fit those molds!

2mo ...

Insightful · 🗣️ 13



Reply · 2 Replies



Peter Mancini (He/Him) · 1st

Principal at Stealth Startup

1. Not knowing what distribution their data is and using tools that depend upon a specific distribution, usually the normal distribution.

2mo ...



Tyler Buffington, PhD · 1st

Senior Data Scientist

1. Not understanding decision analysis and value of information (VOI). For example, building an extremely complicated model when a simple analysis is sufficient to guide the business to the correct decision.
2. Failing to understand that underpowered experiments are untrustworthy even if the result is statistically significant.

2mo ...

Insightful · 🗣️ 6



Reply · 2 Replies



Aaron Sheldon · Following

VP of Engineering at Genstate

In general across all scientific models. Basically applying statistical exploratory analysis. Or more generalizations all the way up to that can tell you how your system

Statistical tests and models assume of how the system your are studying the quantitative bounds of your theory right or wrong. This is even true of cl

Celebrate · 🌍 34

Reply · 4 Replies



Kristen Kehr (She/Her) · 1st

Data Scientist | Developer Advocate

Not so much fellow statisticians "correlation is not causation" stakeholders after sharing

Like · 🗣️ 9

Reply



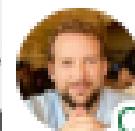
Thomas Speidel · 2nd

Statistician and Data Scientist

1. The independent/predictor variables have to be normally distributed
2. Linear regression can only fit a straight line
3. Logistic regression is a classifier
4. One should always split the data into training and test sets and
5. Cross-validation is the only resampling method
6. Outliers need to be removed
7. One should always remove non-significant predictors
8. R^2 is the only performance metric to look at
9. Doing a vanilla least square regression because that's what we remember from stat 101, and then jumping to a NN with 100 hidden layers and 12 hours of training time because the linear regression didn't perform well
10. That explaining a model simply means doing a feature importance/SHAP plot
12. Not thinking causally
13. Not quantifying uncertainty of the model and predictions
14. Not doing literature reviews (remembering that data science is a vast field that grew in parallel but independent fields)
15. How unstable and sometimes problematic feature selection algos are (at least for weak features)

Like · 🗣️ 13

Reply · 2 Replies



Joe Wyer (He/Him) · 1st

Head of Economics and Applied Science @ Haus | We're Hiring!

Thinking statistical significance should be the filter for go/no go of business decisions.

Like · 🗣️ 4

Reply · 2 Replies

Rodrigo Rivera · 1st

Data Products for Business Users Builder | 🐱 Enabling Anyone t...

around the p-value

ly · 3 Replies



Zion Pibowei • 1st

Principal Data Scientist | Head of Data Science @Periculum | Data ...

Using SMOTE when they should not 🤔

2mo ...

Celebrate

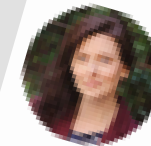


2

Reply

1

2mo ...



Leigh McCormick

Proven

L
do

Insig.

2mo ...

Jeff Bean • 2nd

Senior Consultant at Deloitte

Failing to recognize the difference between "significant" (i.e., $P < 0.05$, not likely due to chance, etc.) and "meaningful" (i.e., Cohen's delta, worthwhile, cost-effective, etc.). You need both.

Celebrate



2

Reply

Replies

life data

2mo ...



Peter Mancini (He/Him) • 1st
Principal at Stealth Startup

1. Not knowing what distribution their data is and using tools that depend upon a specific distribution, usually the normal distribution.



Tyler Buffington, PhD • 1st
Senior Data Scientist

1. Not understanding decision analysis and value of information (VOI). For example, building an extremely complicated model when a simple analysis is sufficient to guide the business to the correct decision.
2. Failing to understand that underpowered experiments are untrustworthy even if the result is statistically significant.

Insightful



6

Reply

2 Replies

2mo ...

Rodrigo Rivera • 1st
Data Products for Business Users Builder | Enabling Anyone t...



6

Reply

2 Replies

around the p-value

3 Replies

...



Aaron Sheldon • Following
VP of Engineering at Genstate

In general across all scientific models. Basically applying statistical exploratory analysis. Or more generalizations all the way up to that can tell you how your system

Statistical tests and models assume of how the system your are studying the quantitative bounds of your theory right or wrong. This is even true of cl

Celebrate



34

Reply

4 Replies



Kristen Kehrer (She/Her) • 1st
Data Scientist | Developer Advocate

Not so much fellow statisticians "correlation is not causation" stakeholders after sharing

Like



9

Reply



Thomas Speidel • 2nd
Statistician and Data Scientist

1. The independent/predictor variables have to be normally distributed
2. Linear regression can only fit a straight line
3. Logistic regression is a classifier
4. One should always split the data into training and test sets and
5. Cross-validation is the only resampling method
6. Outliers need to be removed
7. One should always remove non-significant predictors
8. R^2 is the only performance metric to look at
9. Doing a vanilla least square regression because that's what we remember from stat 101, and then jumping to a NN with 100 hidden layers and 12 hours of training time because the linear regression didn't perform well
10. That explaining a model simply means doing a feature importance/SHAP plot
12. Not thinking causally
13. Not quantifying uncertainty of the model and predictions
14. Not doing literature reviews (remembering that data science is a vast field that grew in parallel but independent fields)
15. How unstable and sometimes problematic feature selection algos are (at least for weak features)

Like



13

Reply

2 Replies

a broken record saying "results" to



Joe Wyer (He/Him) • 1st

Head of Economics and Applied Science @ Haus | We're Hiring!

Thinking statistical significance should be the filter for go/no go of business decisions.

Like




4

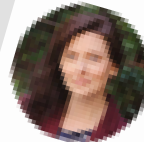
Reply

2 Replies

1mo ...


 **Zion Pibowei** • 1st
Principal Data Scientist | Head of Data Science @Periculum | Data ...
Using SMOTE when they should not 🤔


Celebrate • 🌍 2 | Reply • 1


 **Leigh McCormick**
Proven ...
L ...
Insig.


Jeff Bean • 2nd
Senior Consultant at Deloitte
Failing to recognize the difference between "significant" (i.e., $P < 0.05$, not likely due to chance, etc.) and "meaningful" (i.e., Cohen's delta, worthwhile, cost-effective, etc.). You need both.


Celebrate • 🌍 2 | Reply

 **Jessie Lamontagne** • 1st
Machine Learning | Data Science | Analytics
Describing only the most vanilla ...


 **Thomas Speidel** • 2nd
Statistician and Data Scientist
1. The independent/predictor variables have to be normally distributed
2. Linear regression can only fit a straight line

 **Aaron Sheldon** • Following
VP of Engineering at Genstate
In general across all scienti ...
strong limitations ...

 **Justyn Hornor (He/Him)** • 2nd
Fractional CTO/CPO. Polymath. Inventor. Technology generalist. Str ...
Not recognizing Type 1 vs Type 2 errors as it relates to the org's needs.

 **Peter Mancini (He/Him)** • 1st
Principal at Stealth Startup
1. Not knowing what distribution their data is and using tools that depend upon a specific distribution, usually the normal distribution.


2mo ...

 **Tyler Buffington, PhD** • 1st
Senior Data Scientist
1. Not understanding decision analysis and value of information (VOI). For example, building an extremely complicated model when a simple analysis is sufficient to guide the business to the correct decision.
2. Failing to understand that underpowered experiments are untrustworthy even if the result is statistically significant.


Insightful • 🗣️ 6 | Reply • 2 Replies

... didn't perform well
10. That explaining a model simply means doing ...
12. Not thinking causally importance/SHAP plot
13. Not quantifying uncertainty of the model and predictions
14. Not doing literature reviews (remembering that data science is a vast field that grew in parallel but independent fields)
15. How unstable and sometimes problematic feature selection algos are (at least for weak features)


Like • 🗣️ 13 | Reply • 2 Replies

 **Kristen Kehrer (She/Her)** • 1st
Data Scientist | Developer Advocate
Not so much fellow stat ...
"correlation is not caus ...
stakeholders after shar ...


Like • 🗣️ 9 | Reply

 **Rodrigo Rivera** • 1st
Data Products for Business Users Builder | 🐱 Enabling Anyone t ...
around the p-value

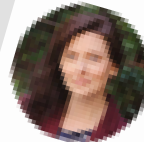
... 3 Replies

 **Joe Wyer (He/Him)** • 1st
Head of Economics and Applied Science @ Haus | We're Hiring!
Thinking statistical significance should be the filter for go/no go of business decisions.

Like • 🗣️ 4 | Reply • 2 Replies


 **Zion Pibowei** • 1st
Principal Data Scientist | Head of Data Science @Periculum | Data ...
Using SMOTE when they should not 🤔


Celebrate • 🌍 2 | Reply • 1

 **Leigh McCormick**
Proven ...
L ...
Insig.

Jeff Bean • 2nd
Senior Consultant at Deloitte
Failing to recognize the difference between "significant" (i.e., $P < 0.05$, not likely due to chance, etc.) and "meaningful" (i.e., Cohen's delta, worthwhile, cost-effective, etc.). You need both.


 **Jessie Lamontagne** • 1st
Machine Learning | Data Science | Analytics
Describing only the most vanilla ...

 **Aaron Sheldon** • Following
VP of Engineering at Genstate
In general across all scienti ...
strong limitations ...


 **Justyn Horner** (He/Him) • 2nd
Fractional CTO/CPO. Polymath. Inventor. Technolo ...
Not recognizing Type 1 vs Type 2 errors as ...
needs.

Like • 🌍 3 | Reply ...
... bounds of your thet ...
right or wrong. This is even true of cl ...

Celebrate • 🌍 34 | Reply • 4 Replie

 **Kristen Kehrer** (She/Her) • 1st
Data Scientist | Developer Advocate
Not so much fellow stat ...
"correlation is not caus ...
stakeholders after sha ...

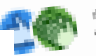
Like • 🌍 9 | Reply

 **Timothy Dobbins** • You
Building data products @ Gridsearch | Principal DS @ Trilliant Health
2mo • Edited • 📝

Statistics is the biggest skills gap in data science right now. I thought it was software engineering but having seen how data science is done in many organizations and hearing from statisticians yesterday about the common mistakes data scientists make, I'm convinced that the coding gaps are nothing compared to the stats gaps.

My stats friends really showed up yesterday.


🌍 🗣️ 🌍 Mary van Valkenburg and 414 others • 63 comments • 31 reposts

 **Peter Mancini** (He/Him) • 1st
Principal at Stealth Startup
1. Not knowing what distribution their data is and using tools that depend upon a specific distribution, usually the normal distribution.


Huffington, PhD • 1st
Data Scientist
Understanding decision analysis and value of information
example, building an extremely complicated model when
analysis is sufficient to guide the business to the correct

... didn't pe ...
10. Tha ...
import ...
12. M ...
13. ...
14. ...
vast field ...
15. How unstable ...
algos are (at least for w ...

Like • 🌍 13 | Reply • 2 Replies

 **Rodrigo Rivera** • 1st
Data Products for Business Users Builder | Enabling Anyone t ...
around the p-value

insightful • 🌍 6 | Reply • 2 Replies

 **Joe Wyer** (He/Him) • 1st
Head of Economics and Applied Science @ Haus | We're Hiring!
Thinking statistical significance should be the filter for go/no go of business decisions.

Like • 🌍 4 | Reply • 2 Replies



REQUIREMENTS ARE HARD TO DEFINE BECAUSE OF **AMBIGUITY**

SUCCESS IS HARD TO DEFINE BECAUSE OF **UNCERTAINTY**

VALUE IS MIS-MEASURED BECAUSE OF **BAD SCIENCE**





AMBIGUITY OCCURS WHEN WE DON'T KNOW WHAT QUESTIONS TO ASK

UNCERTAINTY OCCURS WHEN WE DON'T KNOW HOW TO VALIDATE THE ANSWERS

BAD SCIENCE OCCURS WHEN WE INCORRECTLY VALIDATE THE ANSWERS.



THANK YOU!

LET'S CONNECT

