

## Notes: The Parable of Google Flu — Big Data Pitfalls

Pitfall	Explanation
<b>Overfitting</b>	The model captured noise and coincidental correlations instead of real causal patterns. GFT mistook spikes in “flu symptom” searches as actual outbreaks.
<b>Overparameterization</b>	Too many parameters made the model overly flexible. With thousands of variables and no regularization, GFT fit historical data but failed in prediction.
<b>Lack of Ground Truth Validation</b>	No continuous calibration with CDC data; errors accumulated over time.
<b>Algorithmic Drift</b>	User search behavior and media attention changed, breaking previous correlations.
<b>Correlation ≠ Causation</b>	Search volume increases reflected public concern, not infection rates.
<b>Lack of Transparency</b>	The algorithm and variable list were not publicly released, preventing replication.
<b>Neglecting Traditional Data</b>	Ignoring epidemiological data led to unstable predictions; hybrid models performed better.
<b>Big Data Hubris</b>	The belief that massive data alone can replace theory and domain expertise.