

Data Science Fundamentals Feature Engineering





Sergej Stoppel
in/sergej-stoppel

Today, we're going to look at feature engineering. It's a fancy way of saying we're going to pick out the best ingredients to make the perfect data smoothie. And just like smoothies, you need to choose the right mix of features to create a delicious blend.





Sergej Stoppel
in/sergej-stoppel

What is Feature Engineering?

Feature engineering is the art of selecting, extracting, and transforming features from raw data to create new features that can improve the performance of machine learning models. It's like a chef selecting the best ingredients to make a delicious dish.



Sergej Stoppel
in/sergej-stoppel

Why do we care?

The quality of features used for machine learning models can make or break their performance.

Feature engineering can improve the accuracy of the model, help it generalize better, and reduce overfitting.



Sergej Stoppel
in/sergej-stoppel

Types of Features

There are two main types of features: numerical and categorical. Numerical features are continuous, like height or weight. Categorical features are discrete, like color or gender.

Let's have a look at what we can do to improve the features for our models. The first step is knowing what features you want to use.



Sergej Stoppel
in/sergej-stoppel

First: Feature Selection

Feature selection is the process of selecting the most relevant features for a model. It's like choosing the best toppings for your pizza. You don't want too many, but you also don't want too few.





Sergej Stoppel
in/sergej-stoppel

Sometimes you can't use the data raw. This where feature extraction is your friend.

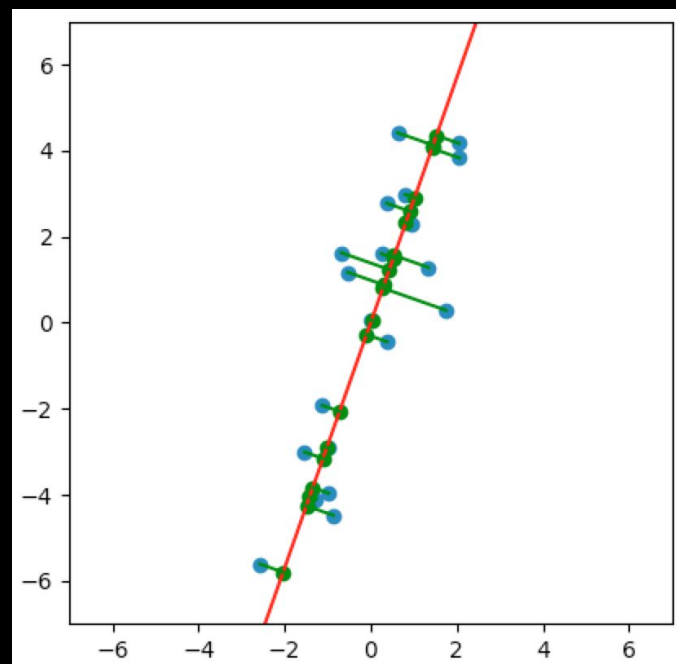
Feature extraction is the process of creating new features from the existing ones. It's like combining and to get a beautiful dish.





Sergej Stoppel
in/sergej-stoppel

Techniques like Principal Component Analysis (PCA) and Singular Value Decomposition (SVD) can help you extract the essence of your data and create new features that capture important patterns.





Sergej Stoppel
in/sergej-stoppel

Another important technique is feature scaling. It is the process of normalizing features so they're all on the same scale.

It's like trying to compare apples to oranges, and realizing that they're both just fruit. Scaling your features can help your model converge faster and reduce the impact of outliers.





Sergej Stoppel
in/sergej-stoppel

Finally we come to feature transformation. It is the process of modifying features to make them suitable for your model.

Personally I see feature scaling as a subset of feature transformation, but you often find the two aspects discussed separately, so I mentioned that as well.





Sergej Stoppel
in/sergej-stoppel

Feature transformation is like taking a square peg and turning it into a round one, so it fits snugly in the hole.

There are soooo many feature transformation techniques, like logarithmic or exponential transformations, binning, mean or max pooling, encoding transformations and so much more.



Sergej Stoppel
in/sergej-stoppel

There are many tools available to help you with feature engineering, from open-source libraries like scikit-learn and TensorFlow, to cloud-based platforms like AWS and Azure.

Whatever your preference, there's a tool out there that can help you build better features faster.



Sergej Stoppel
in/sergej-stoppel

Like any process, feature engineering has its own set of best practices. Here are some tips to keep in mind:

- **Understand your data and the problem you're trying to solve.**
- **Keep it simple: start with basic feature engineering techniques before moving on to more complex ones.**



Sergej Stoppel
in/sergej-stoppel

- **Validate your features: test their impact on your model and iterate as needed.**
- **Don't overfit: avoid creating features that are too specific to your training data, as they may not generalize well to new data.**



Sergej Stoppel
in/sergej-stoppel

Feature engineering is not without its challenges. Here are some common ones:

- **Data quality: feature engineering requires clean, well-structured data. If your data is messy or inconsistent, it can be difficult to create useful features.**
- **Data quantity: feature engineering can require a large amount of data, especially for more complex models.**



Sergej Stoppel
in/sergej-stoppel

- **Time and resources: feature engineering can be time-consuming and require a significant amount of computational resources.**
- **Overfitting: as we mentioned earlier, creating features that are too specific to your training data can lead to overfitting and poor generalization.**



Sergej Stoppel
in/sergej-stoppel

Let's take a look at some real-world examples of feature engineering in action:

- **Fraud detection: By creating clever features that capture sneaky behavior, you can train machine learning models to stop fraud in its tracks. It's like giving your model a detective's intuition, so you can protect your assets and sleep better at night.**



Sergej Stoppel
in/sergej-stoppel

- **Customer segmentation:**
By creating features that capture customer behaviors and demographics, you can train machine learning models to segment your customers for targeted messaging.

It's like a dating app for your customers, so you can boost sales and make them feel like they've found their perfect match.



Sergej Stoppel
in/sergej-stoppel

- **Predictive maintenance:**

By creating clever features that capture performance and usage patterns, you can predict when a machine is likely to fail. It's like superhero vision for your machines, so you can fix them before they break. Get creative with your feature engineering - your budget will thank you!



Sergej Stoppel
[in/sergej-stoppel](#)

Feature engineering is not a static process; it's constantly evolving. With the advent of new techniques like deep learning and automated feature engineering, the future looks bright for this field. Who knows where we'll go next?



Remember

- 1. Feature engineering = making data suited for machine learning models.**
- 2. Different features need different techniques: Scaling, transforming, extracting, you name it.**
- 3. Rules to live by: know your data, keep it simple, test your features, avoid overfitting**

**Feel free to reach out
or to connect with me
for more weekly
slideshows on
visualization, data
science **and** machine
learning.**

Sergej Stoppel
[in/sergej-stoppel](#)

