# Lecture Note 5: Feature Engineering

## 1. Introduction to Feature Engineering

Feature engineering is the process of leveraging domain expertise and analytical techniques to craft and curate features that enhance the performance of data models. The process consists of two primary tasks: feature expansion (creation or generation), which involves creating new features from existing data, and feature selection, which focuses on identifying the most relevant features.

Feature engineering is crucial because it directly impacts the performance of data models. By creating new features, modifying existing ones, and selecting the most relevant ones, feature engineering enhances the ability of algorithms to extract patterns and make accurate predictions from the data. In essence, it optimizes the representation of data, leading to more effective and efficient model outcomes.

Feature engineering is challenging because it requires a deep understanding of both the data and the problem domain. Unlike other aspects of data science that may follow more structured procedures, feature engineering often involves intuition, creativity, and domain expertise to determine which features will be most informative for the model. Additionally, there is no one-size-fits-all approach, as the effectiveness of features can vary depending on the specific problem being addressed and the algorithms being used. Therefore, feature engineering is often described as more of an art than a science, emphasizing the need for experimentation, iteration, and continuous refinement to achieve optimal results.

Feature engineering is often intertwined with other data science steps, particularly data preprocessing. It's not always clearly delineated from these steps and can naturally occur within them. For instance, when preparing a 'date' column for our analysis, we might derive new features like 'yr_sold', 'mon_sold', and 'day_sold'. Similarly, applying one-hot encoding to a categorical variable like 'view' results in creating multiple binary features. Additionally, features such as 'ID' may be ignored or removed during preprocessing. Transformation techniques can also yield new features that are considered integral to the feature engineering process by some practitioners. These examples illustrate how feature engineering integrates with various data preprocessing tasks, highlighting its organic nature within the data science workflow.

This note, however, delves into systematic approaches to feature engineering, emphasizing structured methods and techniques for enhancing the quality and effectiveness of features in machine learning models. The primary focus revolves around two key aspects of feature engineering: feature expansion and feature selection.

## 2. Feature Expansion

Feature expansion involves creating new features from existing ones, often with the aim of enhancing the performance of machine learning models. The purpose of feature expansion is to provide the model with additional information or to represent the data in a more suitable form for the chosen model.

In broad terms, feature expansion can be approached through three main methods: (1) Natural Outcome from Data Preprocessing, (2) Methodical Expansion with Polynomials, Interactions, etc., and

(3) Manual Craft of New Features with Domain Expertise. Below, we'll delve into methods 2 and 3 in detail.

## 2.1 Methodical Expansion

Methodical feature expansion involves deliberate techniques aimed at creating new features to improve model performance. This can include generating polynomial features to capture nonlinear relationships, creating interaction features to represent relationships between variables, or employing techniques like binning to group continuous variables into discrete categories.

Advanced methodical expansion techniques include kernel methods and others, which are beyond the scope of this discussion. These approaches often involve complex mathematical transformations that manipulate the feature space in non-linear ways to capture intricate relationships between variables. While these methods can yield powerful results, they typically require specialized knowledge and computational resources, making them less accessible for general feature engineering tasks.

To illustrate the workings of interaction and polynomial feature expansion, let's consider a straightforward scenario with three original features: 'sqft_living', 'bedrooms', and 'view' from our housing dataset. If we solely focus on interaction terms, the interaction expansion process will generate three additional features: 'sqft_living' * 'bedrooms', 'sqft_living' * 'view', and 'bedrooms' * 'view'. On the other hand, if we only consider second-order polynomial expansion without interaction terms, we'll obtain three additional features: 'sqft_living'^2, 'bedrooms'^2, and 'view'^2.

In the case of third-order polynomial expansion without interactions, we'll have three cubic terms alongside the original features and their squares. When both interactions and polynomial terms are considered, alongside the three interaction terms mentioned earlier, we'll also have additional terms such as 'sqft_living'^2 * 'bedrooms', 'sqft_living'^2 * 'view', 'bedrooms'^2 * 'sqft_living', 'bedrooms'^2 * 'view', 'view'^2 * 'sqft_living', and 'view'^2 * 'bedrooms'.

Polynomial feature expansion, where interaction is considered a specialized case, alongside other methodical feature expansion methods, offer the following advantages:
1. Convenient implementation: One significant advantage of utilizing interaction and polynomial feature expansion lies in its seamless integration with numerous data science software packages, facilitating a straightforward implementation process. These packages often provide built-in functions or libraries specifically designed to handle feature expansion tasks, allowing users to efficiently incorporate interaction terms and polynomial features into their data analysis pipelines.
2. Capturing Non-linear Relationships: Polynomial features enable the modeling of non-linear relationships between variables, allowing for more flexible and nuanced representations of complex data patterns. This capability is particularly valuable when dealing with data that exhibits non-linear behavior, as polynomial features can better capture the underlying relationships compared to linear models.
3. Improved Model Performance: By introducing polynomial terms, the model gains the ability to capture higher-order interactions and non-linear effects that may exist in the data. This can lead to improved model performance, as the model becomes more expressive and better able to fit the training data, potentially reducing bias and variance and enhancing predictive accuracy.
4. Reduced Underfitting: In situations where the relationship between the predictor variables and the target variable is non-linear, using polynomial features can help mitigate underfitting by providing the model with additional flexibility to capture complex patterns in the data. This can

result in a more robust and accurate model that better captures the true underlying data structure.

Methodical feature expansion methods, including polynomial expansions, come with some disadvantages:
1. Increased Dimensionality: Adding polynomial features or interaction terms can significantly increase the dimensionality of the dataset. As the number of features grows, it can lead to the curse of dimensionality, making the dataset more sparse and computationally intensive.
2. Overfitting: With an increase in the number of features, there is a higher risk of overfitting the model to the training data. The model might capture noise or random fluctuations in the training data, leading to poor generalization performance on unseen data.
3. Complexity: The interpretation of models with a large number of polynomial features or interactions becomes more complex. Understanding the relationship between the features and the target variable becomes challenging, making it harder to extract meaningful insights from the model.
4. Computational Cost: Training models with a large number of features requires more computational resources and time. This can be impractical, especially for large datasets or when working with limited computational resources.

## 2.2 Manual Craft of New Features

Manual crafting of new features involves leveraging domain expertise and analytical insights to create meaningful variables that enhance the predictive power of a model. Domain experts possess in-depth knowledge and understanding of the subject matter being studied. They are familiar with the key factors and relationships that drive outcomes in a particular domain. By tapping into this expertise, data scientists can identify relevant variables that may not be directly observable in the dataset but are known to influence the target variable. Analytical insights are derived from exploratory data analysis (EDA) and statistical techniques applied to the dataset. Through EDA, data scientists uncover patterns, trends, and correlations within the data. These insights can reveal potential relationships between variables or highlight areas where additional features could improve model performance.

Based on domain expertise and analytical insights, data scientists manually engineer new features that capture important aspects of the data. This could involve combining existing variables in meaningful ways, transforming variables to better align with the underlying relationships, or introducing domain-specific metrics or indicators. As an example, we engineered a new feature by taking the exponential of the existing feature 'grade'. This decision stems from our observation of the scatter plot, which indicates that the exponential function more accurately captures the relationship between 'grade' and the target variable 'price'.

Manual feature crafting is often an iterative process that involves testing different hypotheses, evaluating the impact of new features on model performance, and refining the feature engineering approach based on feedback. Collaboration between domain experts and data scientists is essential to ensure that the engineered features capture relevant domain knowledge while remaining actionable and interpretable.

In summary, manual crafting of new features with domain expertise and analytical insights involves leveraging subject matter knowledge and data analysis techniques to create meaningful variables that enhance the predictive capabilities of machine learning models.

The manual crafting of new features offers several advantages:
1. Domain Expertise Utilization: This approach allows domain experts to leverage their knowledge and understanding of the data domain to create features that are relevant and meaningful.
2. Flexibility: Manual feature engineering provides flexibility in creating features tailored to specific requirements and hypotheses, enabling the incorporation of nuanced insights that automated methods may overlook.
3. Interpretability: Features crafted manually are often easier to interpret and understand, facilitating the extraction of actionable insights from the data.
4. Control Over Complexity: Manual feature engineering allows for the creation of features that strike a balance between complexity and predictive power, enabling better control over model interpretability and generalization.

However, there are also some disadvantages to manual feature crafting:
1. Subjectivity: The process is subjective and relies heavily on the expertise and intuition of the data scientist, which may introduce biases or overlook important patterns in the data.
2. Time-Consuming: Manual feature engineering can be time-consuming, especially when dealing with large datasets or complex feature transformations, leading to increased development time and resource requirements.
3. Limited Scalability: Manual feature engineering may not scale well to datasets with a large number of features or when dealing with real-time data streams, making it impractical for certain applications or scenarios.

Overall, while manual feature crafting offers the potential for fine-grained control and domain-specific insights, it should be used judiciously in conjunction with automated feature engineering methods to strike a balance between interpretability, predictive performance, and scalability.

# 3. Feature Selection

Feature selection is the process of identifying and selecting a subset of relevant features from the original set of variables in a dataset. The goal of feature selection is to improve the performance of machine learning models by reducing dimensionality, removing irrelevant or redundant features, and focusing on the most informative variables.

Feature selection methods can be broadly categorized into three types:
1. Filter Methods: Filter methods evaluate the relevance of features independently of the predictive model. These methods typically rank features based on statistical measures such as correlation, mutual information, or significance tests. Features are then selected or removed based on predefined criteria, such as a specified threshold or a fixed number of top-ranked features.

2. Wrapper Methods: Wrapper methods evaluate subsets of features by training and testing predictive models using different combinations of variables. These methods use the performance of the predictive model, such as accuracy or error rate, as the criterion for selecting features. Common wrapper methods include forward selection, backward elimination, and recursive feature elimination.

3. Embedded Methods: Embedded methods integrate feature selection directly into the model training process. These methods incorporate feature selection as part of the model fitting

algorithm, automatically selecting the most relevant features during training. Examples of embedded methods include Lasso regression, decision trees, and random forests.

Feature selection offers several benefits, including:
1. Improved Model Performance: By removing irrelevant or redundant features, feature selection can improve the accuracy, efficiency, and interpretability of machine learning models.
2. Reduced Overfitting: Feature selection helps mitigate the risk of overfitting by focusing on the most informative features and reducing the complexity of the model.
3. Faster Training and Inference: By reducing the dimensionality of the data, feature selection can lead to faster model training and inference, particularly for high-dimensional datasets.
4. Enhanced Interpretability: By selecting a subset of relevant features, feature selection can simplify the model and make it easier to interpret, understand, and communicate the relationships between input variables and the target variable.

However, feature selection also has some limitations, including the potential for information loss if important features are removed, increased computational complexity for certain methods, and the need for careful validation to ensure that the selected features generalize well to unseen data. Overall, feature selection is a crucial step in the machine learning pipeline, helping to improve model performance and facilitate better understanding of the underlying data patterns.

# 4. EDA and Feature Engineering

Exploratory Data Analysis (EDA) plays a pivotal role in facilitating feature engineering, encompassing both feature expansion and feature selection processes.

1. Feature Expansion: EDA aids in identifying patterns, trends, and relationships within the data, which can inform the creation of new features. By exploring the relationships between variables, EDA can unveil potential interactions or nonlinearities that suggest avenues for feature expansion. For instance, scatter plots, correlation matrices, and heatmap visualizations can reveal correlations between variables, guiding the creation of interaction terms or polynomial features. Moreover, EDA techniques such as box plots and histograms can uncover distributions and outliers, prompting transformations or the creation of categorical variables based on specific thresholds. Overall, EDA provides the groundwork for generating new features that capture the underlying complexity of the data.

2. Feature Selection: EDA assists in identifying the most relevant features for predictive modeling by uncovering their relationships with the target variable. Through visualizations and statistical analyses, EDA helps assess the significance of individual features in explaining the variation in the target variable. Techniques like correlation analysis, hypothesis testing, and visualization of feature distributions across different target variable classes aid in understanding feature importance. EDA can also reveal multicollinearity among features, guiding the selection of a subset of predictors that offer unique information. By iteratively exploring the data and evaluating feature relevance, EDA facilitates the selection of a concise yet informative feature set, enhancing model interpretability and performance.

In our exploration of Exploratory Data Analysis (EDA), we underscored the significance of employing various techniques to capture relationships between pairs of variables. Among these techniques, two prominently featured methods are correlation coefficient analysis and hypothesis testing, encompassing tests of independence and t-tests.

Correlation coefficient analysis serves as a foundational pillar in EDA, enabling us to quantify the strength and direction of the linear relationship between two continuous variables. By computing correlation coefficients such as Pearson's correlation coefficient and Spearman's rank correlation, we gain valuable insights into the degree of association between variables. This assessment aids in identifying potential predictor variables that exhibit significant correlations with the target variable, thereby guiding feature engineering efforts.

Hypothesis testing, on the other hand, provides a rigorous statistical framework for evaluating the relationship between categorical variables. Through tests of independence, such as the Chi-Square test, we ascertain whether there exists a statistically significant association between two categorical variables. Similarly, the t-test enables us to compare the means of two groups and determine whether the observed differences are statistically significant. These hypothesis tests play a pivotal role in discerning the relevance of individual features within the dataset, facilitating the identification and selection of informative variables for predictive modeling tasks.

By integrating correlation coefficient analysis and hypothesis testing into our EDA toolkit, we augment our ability to uncover meaningful relationships and patterns within the data. These methods serve as essential components of feature engineering, aiding in the creation and identification of relevant features that contribute to predictive model performance, and ultimately enhancing the efficiency and efficacy of our machine learning endeavors.