# Report on end to end machine learning pipeline

**End-to-End Data Science Project (Example: Real Estate Prices)**

1. Understand the Problem (Define Objective):
   Goal: Predict house prices based on features like size, bedrooms, and location.
   Example: A client wants to estimate the price of a new housing.

2. Raw Data Source → Data Ingestion (Collect Data):
   Gather historical real estate data from company records or public sources.
   Example: CSV file with house features and sale prices.

3. Data Evaluation (Explore & Visualize Data):
   Analyze data quality, check distributions, correlations, and outliers.
   Example: Plot size vs. price to observe trends

4. Prepare Data (Data Cleaning & Preprocessing):
   Clean missing values, encode categorical variables, scale and normalize features.
   Example: Replace missing bedroom counts with median values.

5. Modeling (Select & Train Model):
   Choose a model like Linear Regression and train on the dataset.
   Example: Fit a regression model to learn the relationship between house features and prices.

6. Evaluation (Fine-Tune & Validate):
   Assess model performance using metrics (RMSE, $R^2$) and optimize hyperparameters.
   Example: Use cross-validation to select the best features or adjust regularization to reduce prediction error.

7. Deployment (Present Solution & Launch Model):
   Share predictions and insights in a report/dashboard and deploy the model for real-world use.
   Example: A dashboard highlighting top factors influencing prices; integrate the model into a web app where clients input house details to get price predictions.

8. Monitoring & Maintenance:
   Continuously monitor model performance, update new data, and retrain when accuracy drops.
   Example: Track prediction errors over time; retrain the model timely with new data.

Data Pipelines in Machine Learning:

A data pipeline is a sequence of data processing steps, commonly used in machine learning systems to handle large amounts of data and multiple transformations.

How it works:
- Each component in the pipeline works independently and asynchronously.
- A component pulls data from a data store, processes it, and outputs the result to another data store.
- The next component later takes this output as input and continues the process.

Advantages:
- Components are self-contained, making the system easier to understand and maintain.
- Different teams can focus on different components.
- If one component fails, downstream components can often continue using the last available data, improving system robustness.

Challenges:
- If proper monitoring is not in place, a broken component may go unnoticed.
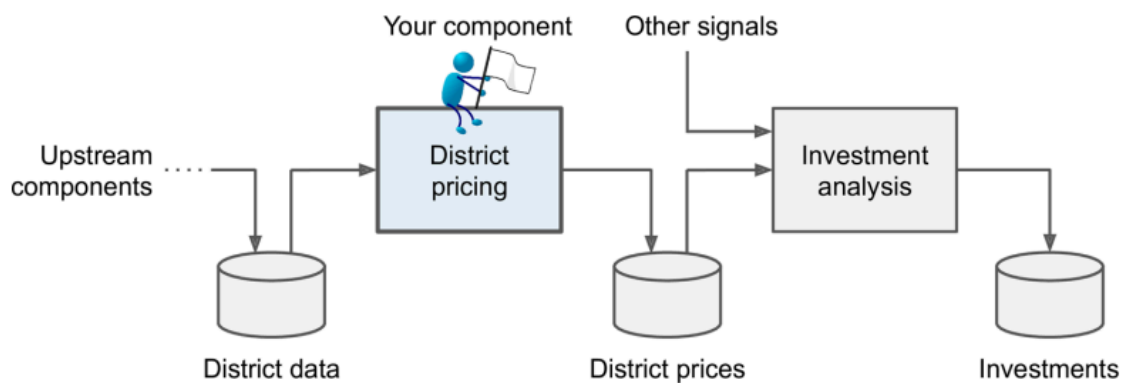- Stale data can reduce overall system performance.



Figure 2-2. A Machine Learning pipeline for real estate investments

When learning Machine Learning, it is more effective to practice with real-world datasets rather than artificial ones. Fortunately, there are thousands of open datasets available across different domains. Some of the most popular sources include the UC Irvine Machine Learning Repository, Kaggle datasets, and Amazon's AWS open data. You can also explore meta portals that list multiple repositories, as well as dedicated data portals such as OpenDataMonitor and Quandl. Additionally, there are curated lists on Wikipedia, community discussions on Quora, and user-shared datasets on Reddit's datasets. These platforms provide a wide range of options, making it easier to experiment with real data and build practical ML skills.

# End to end ml pipeline steps in detail:

When starting a Machine Learning project, the first step is to **clarify the business objective**, since building a model is not the end goal but a means to create value. E.g. in the case of calculating the district's housing prices, the model will predict a district's median housing price, and its output will feed into another system that decides whether to invest in that area, directly impacting revenue. Currently, experts manually estimate housing prices using complex rules, which is costly, slow, and often inaccurate by over 20%. To improve this, census data with housing prices and related features will be used to train a model. Framing the problem, it is a supervised learning regression task (more specifically, a multiple regression predicting one target value), and since the dataset fits in memory with no need for real-time updates, batch learning is sufficient.

Since this project is about predicting numerical values (median housing prices), we need a **performance measure** to check how good the model's predictions are. For regression problems, a very common choice is Root Mean Square Error (RMSE). RMSE tells us how far, on average, our predictions are from the actual values. Importantly, it penalizes large errors more heavily than small ones, which is useful here because a big mistake in predicting house prices could lead to bad investment decisions.

$$\text{RMSE}\left(\mathbf{X}, h\right) = \sqrt{\frac{1}{m} \sum_{i=1}^{m} \left(h\left(\mathbf{x}^{(i)}\right) - y^{(i)}\right)^2}$$

It's also important to **check assumptions**—for example, whether the downstream system needs exact prices or just categories. After confirming they need actual prices, the regression approach is validated, and development can begin.

In most real projects, data is stored in databases or spread across multiple files, and you would need access permissions and knowledge of the schema. For this housing price project, however, it's much simpler: the data comes in a single compressed file (housing.tgz) containing a CSV file (housing.csv). While you could manually **download and extract** it, it's better to create a function that automates the download and extraction. This way, if the data updates regularly or needs to be installed on multiple machines, the process can run automatically and consistently.

Even early in a project, it's important to **set aside a test set** to evaluate your model later. If you look at the test data while exploring or selecting models, you risk data snooping bias, which makes performance estimates too optimistic. A common approach is to randomly select about 20% of the data as the test set, but simple random splits can change each run. To make it stable,

you can save the test set, set a random seed, or use a unique identifier for each instance (e.g., a hash of the row ID or stable features like latitude/longitude). This ensures consistency even if new data is added. Libraries like Scikit-Learn provide functions such as train_test_split() to simplify this process, allowing you to split multiple datasets consistently and set the random seed for reproducibility.

At this stage in the pipeline, the goal is to **explore and visualize the data** to gain deeper insights. Make sure the test set is set aside and only the training set is used for exploration. If the dataset is very large, you can work on a smaller sample to make analysis faster, but for small datasets, you can use the full training set. It's a good practice to create a copy of the training set so you can experiment and manipulate the data without affecting the original. This helps you safely analyze distributions, correlations, and patterns before modeling.

To understand how features relate to each other and to the target variable, you can **compute correlations** between all pairs of attributes using the Pearson correlation coefficient (r), which ranges from −1 to 1. A value close to 1 indicates a strong positive correlation (e.g., median house value rises with median income), a value close to −1 indicates a strong negative correlation (e.g., prices slightly decrease as latitude increases), and values near 0 mean little or no linear correlation. Visualizing these relationships with plots alongside their correlation coefficients helps quickly identify which features are most related to the target and which may be less useful for modeling.

After exploring correlations and distributions, the next step is **experimenting with attribute combinations** to create more informative features. Some original attributes, like the total number of rooms or bedrooms, may not be very useful by themselves, but combining them with other attributes can reveal stronger relationships with the target. For example, computing bedrooms per room, rooms per household, or population per household can provide better predictors of median house value. In fact, the new bedrooms_per_room attribute often shows a stronger correlation with housing prices than the raw counts, indicating that houses with lower bedroom-to-room ratios tend to be more expensive. This exploratory step helps you gain insights and prepare better features for your initial model, and it's an iterative process—you can refine these combinations as you analyze the prototype's output.

The next step is to **prepare the data** for Machine Learning algorithms. Instead of transforming the data manually, it's best to write functions for this purpose. This approach allows you to reproduce transformations easily on new datasets, build a reusable library of functions for future projects, and apply the same transformations to data in a live system. It also makes it easier to experiment with different transformations to see which combination works best. To start, revert to a clean copy of the training set and separate the predictors from the labels, since the same transformations should not always be applied to both.

Most Machine Learning algorithms cannot handle missing values, so we need to fix them by applying **data cleaning** techniques. For attributes like total_bedrooms, there are three options: drop the districts, drop the attribute, or fill missing values with a value such as zero, mean, or median. The recommended approach is to compute the median on the training set and use it to fill missing values consistently, saving the median for later use on the test set or new data. Scikit-Learn's SimpleImputer simplifies this: create an instance with strategy="median", fit it to the numerical features (excluding text attributes), and then transform the dataset. This replaces missing values with the learned medians. The resulting array can be converted back into a pandas DataFrame for further processing.

**Scikit-Learn** is ideal for end-to-end ML pipelines because it provides a consistent, simple API for all steps: data preprocessing, feature transformation, model training, prediction, and evaluation. Its estimators, transformers, and pipelines make it easy to chain steps, handle missing values, and automate workflows. With sensible defaults, composability, and reproducibility, you can quickly build, test, and deploy complete ML pipelines efficiently.

Its main principles include:

Consistency: All objects follow a simple, uniform interface.
- Estimators: Objects that learn parameters from data via fit(). Hyperparameters are set when creating the instance.
- Transformers: Estimators that can modify data using transform(); fit_transform() combines fitting and transforming efficiently.
- Predictors: Estimators that make predictions with predict() and can evaluate performance with score().
- Inspection: Hyperparameters and learned parameters are accessible via public variables (learned parameters often end with _).
- Nonproliferation of classes: Data is handled as NumPy arrays or SciPy matrices; hyperparameters are standard Python types.
- Composition: Building blocks can be combined, e.g., using Pipeline to chain transformers and a final estimator.
- Sensible defaults: Most parameters have reasonable defaults, allowing quick creation of baseline models.

In context of **text and categorical attributes**, some features, like ocean_proximity, are categorical rather than numerical. Most ML algorithms require numbers, so these categories must be converted. One approach is ordinal encoding, which assigns a unique number to each category, but this can mislead algorithms into assuming an order that doesn't exist. A better approach is one-hot encoding, which creates a separate binary feature for each category, only one

is "hot" (1) per instance. Scikit-Learn provides the OneHotEncoder for this, which outputs a sparse matrix to save memory when there are many categories.

For categorical attributes with a large number of categories, one-hot encoding can create very high-dimensional input. Alternatives include converting categories into meaningful numerical features (e.g., distance to the ocean) or using learnable embeddings, where each category is represented by a low-dimensional vector learned during training.

Sometimes you need **custom transformations** that Scikit-Learn doesn't provide, such as combining attributes or performing specific cleanup steps. To make a custom transformer compatible with Scikit-Learn (and pipelines), create a class implementing fit(), transform(), and optionally fit_transform(). You can inherit from TransformerMixin to get fit_transform() for free, and from BaseEstimator to get get_params() and set_params() for hyperparameter tuning.

For example, a transformer can add a new feature like bedrooms_per_room, controlled by a hyperparameter add_bedrooms_per_room=True. This makes it easy to test whether adding the attribute improves model performance. Automating these steps allows you to efficiently experiment with many feature combinations.

**Feature scaling** is crucial because most ML algorithms perform poorly when numerical attributes have very different ranges. For example, in the housing dataset, total rooms and median incomes may have an incomparable range which needs feature scaling. Two common scaling methods are:

1. Min-max scaling (normalization): Rescales values to a fixed range, usually 0–1, using MinMaxScaler in Scikit-Learn.

2. Standardization: Subtracts the mean and divides by the standard deviation, resulting in zero mean and unit variance using StandardScaler. Standardization is less sensitive to outliers, while min-max scaling can be heavily affected by extreme values.

Scaling is typically applied to input features, not the target values.

**Data preprocessing** often involves multiple steps that must be executed in order. Scikit-Learn's Pipeline class helps chain transformers and the final estimator into a single workflow. Each step (except the last) must be a transformer with a fit_transform() method, and calling fit() on the pipeline applies each step sequentially.

For datasets with both numerical and categorical columns, the ColumnTransformer allows applying different transformations to each column type in one step. For example, numerical

columns can go through a pipeline of imputation and scaling, while categorical columns can be one-hot encoded. The outputs are concatenated into a single matrix, which can be dense or sparse depending on the data. This setup creates a complete preprocessing pipeline that can handle the entire dataset automatically, making it easy to reuse, maintain, and integrate with ML models.

For older Scikit-Learn versions (<0.20), alternatives include FeatureUnion or third-party libraries like sklearn-pandas, though they require workarounds to handle different column types.

After preparing the data, you can select and **train ML models**. Start with a simple model like Linear Regression; it works but may underfit, producing high errors. A more complex model like DecisionTreeRegressor can capture nonlinear relationships, but it often overfits the training data. To evaluate models reliably, use K-fold cross-validation, which splits the training set into folds and gives a better estimate of performance and its variability. Comparing models shows that Decision Trees overfit while Linear Regression underfits. A RandomForestRegressor, an ensemble of many trees, often performs better and reduces overfitting. To improve further, you can simplify or regularize models, add more data, or try other algorithms like SVMs or neural networks. The goal is to shortlist a few promising models before fine-tuning.

Once you have a shortlist of promising models, the next step is **hyperparameter tuning**. GridSearchCV tests all combinations of specified hyperparameter values using cross-validation, while RandomizedSearchCV samples a fixed number of random combinations, which is faster for large search spaces. Another approach is ensemble methods, combining top-performing models to improve accuracy.

After fine-tuning, analyze errors, adjust features if needed, and finally **evaluate the model** on the test set using the full pipeline (without refitting). You can also compute a confidence interval to estimate the precision of the model's performance. Avoid tuning based on the test set, as this can lead to overfitting. Before deployment, document the solution, assumptions, limitations, and insights (e.g., which features matter most) to communicate results clearly and support decision-making.

Once approved, **deploy your trained model** (including preprocessing pipelines) to production, either within your application or via a web service/REST API. Cloud platforms like Google Cloud AI Platform can simplify scaling and load balancing.

Deployment is not the end: continuously **monitor the model's performance** to detect sudden drops or gradual decay ("model drift"), which can occur as the data or environment changes over time. Monitoring can rely on downstream metrics, human raters, or automated checks of input data quality. Automate retraining, evaluation, and deployment to handle evolving data, and maintain backups of models and datasets to allow quick rollbacks if needed.

By following these steps systematically, we can build an **end-to-end ML pipeline** that is reproducible, robust, and ready for production. Each step ensures that your data and model are properly prepared, evaluated, and maintained, reducing errors and improving performance over time.

**Mermaid diagram:**