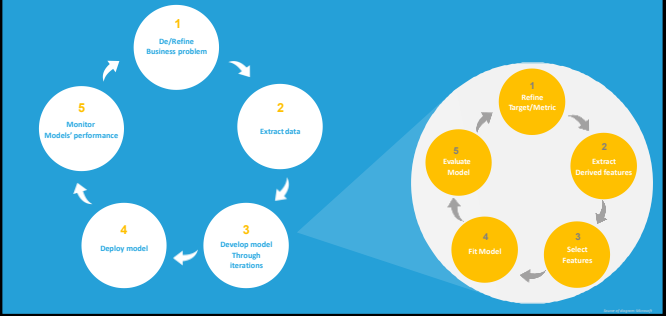


getting started with
Machine Learning

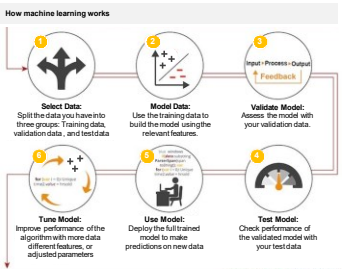
:: Machines Learning uses historical data to make predictions



:: Steps to build a machine learning solution



:: How machine learning works



:: Data in the real world is dirty

Data is rarely clean and often you can have data quality issues

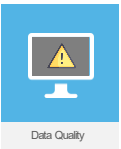
#	Id	Name	Birthday	Gender	IsTeacher?	#students	Country	City
1	111	John	15/12/1995	M	1	10	Canada	London
2	222	Peter	15/10/1978	F	1	15	Canada	London
3	333	Maria	18/04/2000	F	0	0	Spain	Madrid
4	444	John	12/11/1981	M	0	0	Spain	Madrid
5	555	John	15/03/2000	F	1	25	Germany	Berlin
6	666	John	18/03/1975	M	1	10	Italy	Rome
7	777	Carin	25/05/1999	M	0	0	Italy	Rome
8	888	Rebecca	03/06/1945	F	0	0	Portugal	Lisbon
9	999	John	05/06/1983	F	0	5	Netherlands	Amsterdam
10	10/10	Paul	14/11/1992	M	1	25	Italy	Rome

Uniqueness Formats Attribute dependencies Missing values Invalid values Misfielded values Misspellings

- The typical data quality issues that arise are:
- **Incomplete:** Data lacks attributes or containing missing values.
 - **Noisy:** Data contains erroneous records or outliers.
 - **Inconsistent:** Data contains conflicting records or discrepancies.

:: Data in the real world is dirty

What kind of issues affect the quality of data?



- **Invalid values**
Some datasets have well-known values, e.g. gender must only have "F" (Female) and "M" (Male). In this case it's easy to detect wrong values.
- **Formats**
The most common issue. It's possible to get values in different formats like a name written as "Name, Surname" or "Surname, Name".
- **Attribute dependencies**
When the value of a feature depends on the value of another feature. For example, if we have some school data, the "number of students" is related to whether the person "is teacher". If someone is not a teacher teacher can't have any students.
- **Uniqueness**
It's possible to find repeated data in features that only allow unique values. For example, we can't have two products with the same identifier.
- **Missing values**
Some features in the dataset may have blank or null values.
- **Misspellings**
Incorrectly written values.
- **Misfielded values**
When a feature contains the values of another.

:: Preprocessing data - Clean your data

Why to deal with missing values?

- Missing values in a dataset can be due to error or because observations that were not recorded
- When missing value are present, certain algorithms may not work or you may not have the desired result.
- Missing data affects some models more than others¹
- Even for models that handle missing data, they can be sensitive to it (missing data for certain variables can result in poor predictions)²



Missing value is probably the most common problems in data mining/machine learning

:: Preprocessing data - Clean your data

How to deal with missing values?

Typical missing value handling methods are:

- **Deletion**
Remove records with missing values
- **Dummy substitution**
Replace missing values with a dummy value, e.g. unknown for categorical or 0 for numerical values.
- **Mean substitution**
If the missing data is numerical, replace the missing values with the mean.
- **Frequent substitution**
If the missing data is categorical, replace the missing values with the most frequent item
- **Regression substitution**
Use a regression method to replace missing values with regressed values.



:: Preprocessing data - Clean your data

What you should know about the outliers/anomalies

- Outliers may bring about problems by distorting the predictive model.
- What's an outlier is somewhat subjective.¹
- Outliers can be very common in multidimensional data.²
- Some models are less sensitive (more robust) to outliers than others.³
- Outliers can be result of bad data collection, or they can legitimate extreme (or unusual) values.⁴
- Sometimes outliers are the interesting data points we want to model, and other times they just get in the way.⁵



An outlier is a data point that distinctly separate from the rest of the data.

:: Preprocessing data - Clean your data

Cause of outliers



Data from different classes



Data measurement and collection Errors



Natural variation

:: Preprocessing data - Clean your data

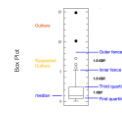
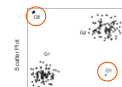
How to deal with outliers?

The choice of how to deal with an outlier should depend on the cause.

- **Keep outliers**
Outliers should not necessarily be omitted from the analysis as they may be genuine observations in the data.

In many applications, outliers provide crucial information. For example, in a credit card fraud detection app, they indicated purchases that fall outside a customer's usual buying patterns.¹
- **Exclude outliers**
There are two common approaches to exclude outliers²
 - **Trimming/Truncation:** Trimming discards the outliers
 - **Winsorizing:** Winsorizing replaces the outliers with the nearest "non-suspect" data

Popular plots for outlier detection:
Scatter Plot / Box Plot



:: Preprocessing data - Clean your data

Examples on how to deal with outliers

• Example for Trimming

Eliminate the outliers "2" & "22"



• Example for Winsorizing

Assign outlier the next highest or lowest value found in the sample that is not an outlier. In this example, "10" & "14" are not outliers and used to replace the outliers "2" & "22".



• Trimming or Winsorizing less than 5% of data points

It will not likely affect the hypothesis testing outcome.

• Trimming or Winsorizing great than 5% of data points

Trimming or Winsorizing great than 5% may affect the outcome results.

- Reduce the power of analysis
- Makes samples less representative
- May affect normality of data
- Consider transforming data, choosing an alternate outcome variable or data analysis technique.

:: Preprocessing data - Data normalization

How to normalize data?

Data normalization re-scales numerical values to a specified range. Popular data normalization methods include:

• Min-Max Normalization:

Linearly transform the data to a range, say between 0 and 1, where the min value is scaled to 0 and max value to 1.

$$X_{\text{norm}} = \frac{X - X_{\text{min}}}{X_{\text{max}} - X_{\text{min}}}$$

• Z-score Normalization (or Standardization):

Scale data based on mean and standard deviation: divide the difference between the data and the mean by the standard deviation. Feature standardization makes the values of each feature in the data have zero-mean (when subtracting the mean in the numerator) and unit-variance: 1.

$$z = \frac{X - \mu}{\sigma} \quad \left\{ \begin{array}{l} \text{Z-Score} = \frac{\text{Value} - \text{mean}}{\text{Standard Deviation}} \end{array} \right.$$

• Decimal scaling:

Scale the data by moving the decimal point of the attribute value.

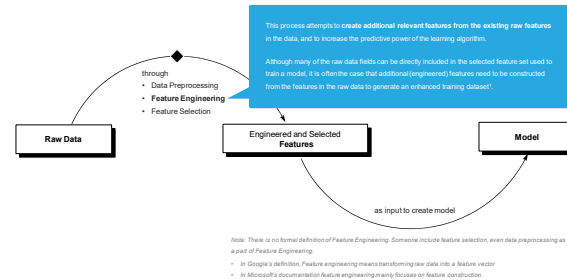
Age	Salary
23	15000
30	20000
35	25000
40	30000
45	35000
50	40000
55	45000
60	50000
65	55000
70	60000

Value Range: 15 - 95 Value Range: \$1,000 - \$15,000

Different scales

Different scales in a dataset may be problematic. In some cases where certain machine algorithms require data to be in the same scale.

:: Feature Engineering can augment your data



:: What's feature ?

Features (input)
Also known as Attributes - Explanatory Variables - Independent Variables, Predictors

Target (output)
Also known as Dependent Variables

	Bedrooms	Sq. feet	Neighborhood	Sales price
Training Data	3	2000	Normalton	\$250,000
	2	800	Hopedorn	\$300,000
	2	850	Normalton	\$150,000
	1	550	Normalton	\$78,000
	4	2000	Skid Row	\$150,000
Test Data	3	2000	Hopedorn	\$220,000
	3	900	Hopedorn	\$800,000
New Data	4	2100	Normalton	\$270,000
	2	450	Hopedorn	
	4	1950	Normalton	

untabeled examples

This unknown value should be predicted by the model

a feature is an individual measurable property or characteristic of a phenomenon being observed*.

In this example, there are a couple of features like bedrooms, house size, neighborhood which are used as input for the modeling. The target is the sales price of the house.

There are some known price for some houses (Training data & Test Data).

A model can be built to predict the the sales price of the houses which are still unknown.

In a nutshell, once we've trained our model with labeled examples, we use that model to predict the label on unlabeled examples.*

:: Feature Engineering - example

Example 2

Raw Data

Additional features are created through feature engineering

• Date features: year, month, week of month, etc.

• Time features

• Season features

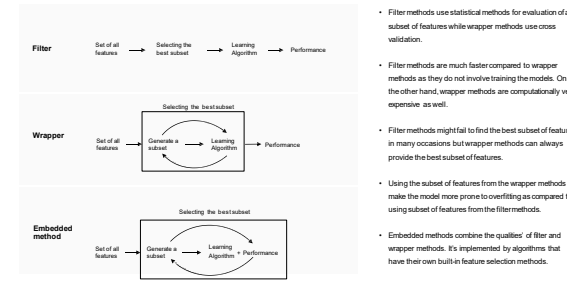
• Weekday and weekend features

• Holiday features: New Year, US Labor Day, US Thanksgiving, Cyber Monday, Christmas, etc.

• Proximal features to capture seasonality

The template provides two ID fields, which also associates with the store ID and productSKU ID.

:: Summary: the modern approaches for Feature Selection



:: Feature Engineering - example

The comparison of the performance results of the four models which are based on 4 training sets

Features	Mean Absolute Error	Root Mean Squared Error	Coefficient of Determination
Feature Set A (training set1) baseline = weekday + holiday + weekday + weekend features for the predicted day	89.7	124.9	0.6
Feature Set A+B (training set2) baseline = previous 12 hours demand	51.7	86.3	0.8
Feature Set A+B+C (training set3) baseline = previous 12 hours demand + previous 12 days at the same hour	47.6	81.1	0.8
Feature Set A+B+C+D (training set4) baseline = previous 12 hours demand + previous 12 days at the same hour + previous 12 weeks at the same hour and the same day demand	46.3	82.1	0.8

Same algorithm but different performance

we used the Boosted Decision Tree Regression module, a commonly used nonlinear algorithm, to build the models.

To understand the performance of four models, see the comparison results in the following table.

- The best results seems from the combination of features A+B+C and A+B+C+D.
- Feature set D does not provide additional improvement over A+B+C.

:: Feature Selection vs Dimension Reduction

Feature Selection

Dimension Reduction

Although feature selection does seek to reduce the number of features in the dataset used to train the model, it is not usually referred to by the term "dimensionality reduction".

Feature selection methods extract a subset of original features in the data without changing them.

Dimensionality reduction methods employ engineered features that can transform the original features and thus modify them.

Examples of dimensionality reduction methods include Principal Component Analysis, canonical correlation analysis, and Singular Value Decomposition.

Modelling

:: Train the model

The modelling is an iterative process of building model which is a combination of data structure, algorithm, and mathematics to captures the relationship between features and target within the dataset

:: What's a model ?

A model defines the relationship between features and label (i.e. target)

Machine learning uses a model to capture the relationship between **feature vectors** and some **target variables** within a training data set.

A feature vector is a set of features or attributes that characterize a particular object, such as the number of bedrooms, bathrooms, and location of an apartment. The target is either a scalar value like rent price, or it's an integer classification such as "creditworthy" or "it's not cancer."

Model

Examples

- An machine learning model is a mathematical model that generates predictions by finding patterns in your data.
- At a high level, a model is a simplification of something more complex.
- A machine learning algorithm use data to automatically learn the rules. It simplifies the complexity of the data into relationships described by rules.
- Modelling is part art and part science.

"All models are wrong, but some are useful."

Machine Learning is to learn a project function for the task

- Speech Recognition
 $f(\text{audio waveform}) = \text{"How are you"}$
- Image Recognition
 $f(\text{cat image}) = \text{"Cat"}$
- Playing Go
 $f(\text{Go board state}) = \text{"5-5" (next move)}$
- Dialogue System
 $f(\text{"How are you?" (what the user said)}) = \text{"I am fine." (system response)}$

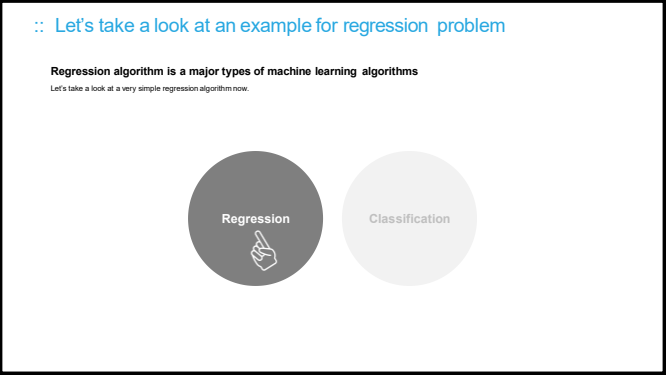
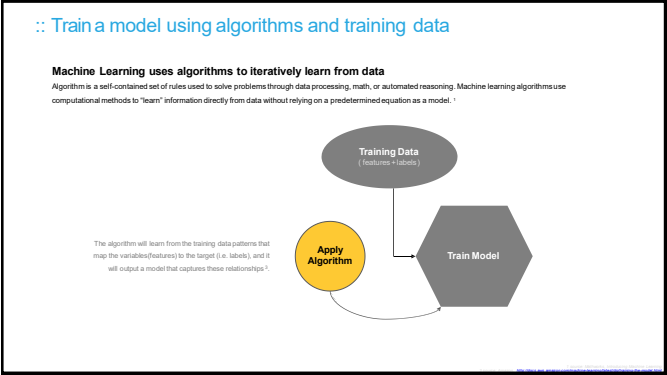
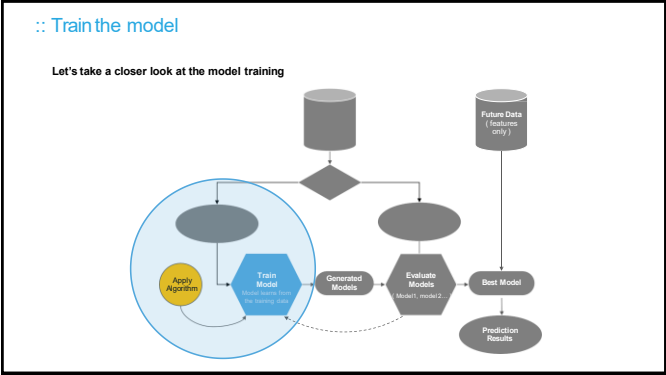
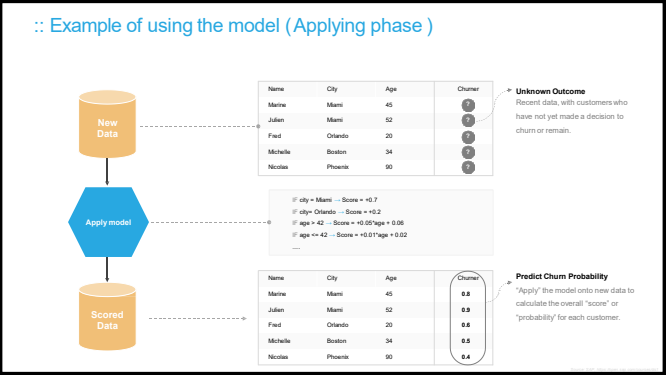
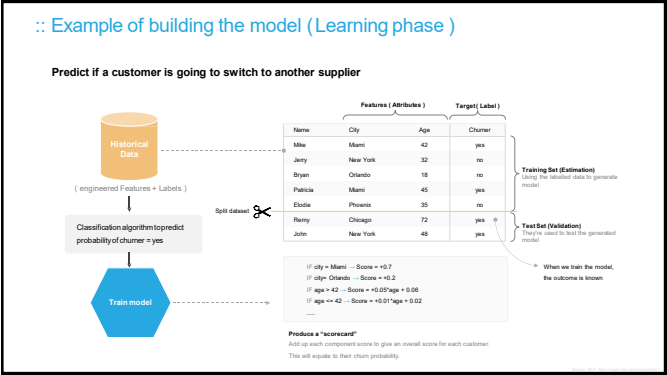
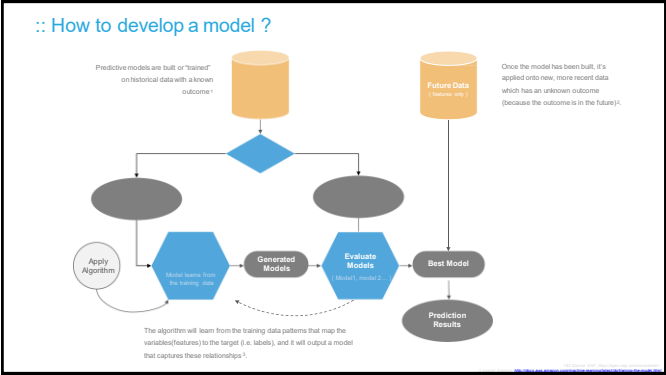
:: The process of Machine Learning

Iterate until data is ready

Iterate to find the best model

Train model, Evaluate Model & Optimization

- Find the model that answers the question most accurately by comparing their success metrics
- Determine if your model is suitable for production



:: Let's take a look at an example for regression problem

Example: Predict house's price using linear regression

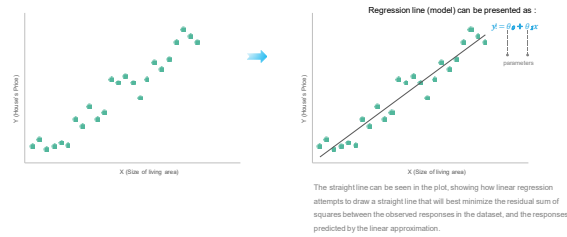
Suppose we have a dataset giving the living areas and prices of 21,613 houses from House Sales in King County, USA. Given data like this, we can learn to predict the prices of other houses in King County.



:: Example: Train a model using Linear regression algorithm

Use linear regression algorithm to approximate the relationship between x and y

Take linear regression as example, the algorithm is trying to find a best-fit line to represent the relationship between the input feature x and target y.



:: Sample python code for building a very simple model

```
import matplotlib.pyplot as plt
import numpy as np
from sklearn import datasets, linear_model

# Load the diabetes dataset
diabetes = datasets.load_diabetes()

# Use only one feature
diabetes_X = diabetes.data[:, np.newaxis, 2]

# Split the data into training/testing sets
diabetes_X_train = diabetes_X[:-20]
diabetes_X_test = diabetes_X[-20:]

# Split the targets into training/testing sets
diabetes_y_train = diabetes.target[:-20]
diabetes_y_test = diabetes.target[-20:]

# Create linear regression object
reg = linear_model.LinearRegression()

# Train the model using the training sets
reg.fit(diabetes_X_train, diabetes_y_train)

# The coefficients
print('Coefficients: \n', reg.coef_)
# The mean squared error
print('Mean squared error: %.2f'
      % np.mean((reg.predict(diabetes_X_test) - diabetes_y_test) ** 2))
# Equivalent non-linear score: R^2 is perfect prediction
print('Score: %.2f' % reg.score(diabetes_X_test, diabetes_y_test))

# Plot results
plt.scatter(diabetes_X_test, diabetes_y_test, color='black')
plt.plot(diabetes_X_test, reg.predict(diabetes_X_test), color='blue',
        linestyle='--')

plt.axis([0, 400, 0, 3000])
plt.show()
```

Important:

Developing a model is a process of experimentation and incremental adjustment as much as it is of applying algorithms to a problem.

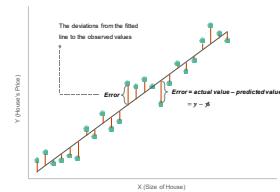
You should expect to spend a lot of time refining and modifying your models to get the best results.

It's very important that you establish a threshold of success for your model before you begin because otherwise you may not know when to stop refining.

:: Usually there is deviation between the actual value and prediction

The deviations indicate how bad the model's prediction was on the training examples

Loss (i.e. error) is a number indicating how bad the model's prediction was on a single example. If the model's prediction is perfect, the loss is zero; otherwise, the loss is greater.



Mean square error (MSE) is a commonly-used function to measure how large the loss is. It's called as **Loss function** or **Cost function**.

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

\hat{y}_i is the prediction
 y_i is the actual value

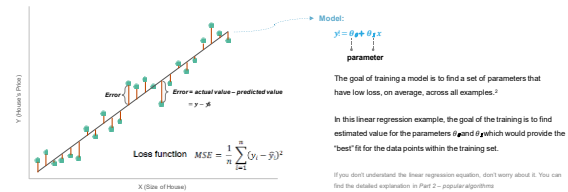
Mean square error (MSE) is the average squared loss per example over the whole dataset. We will elaborate the details of loss function later on.

The smaller the Mean square error, the better the fit of the line to the data.

:: Train the model to minimize the loss/error

Training the model is an iterative process of finding the "best" parameters to minimize the error

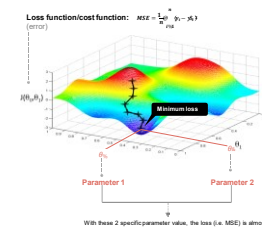
Training a model simply means learning (determining) good values for all the parameters of the model from labeled examples. In supervised learning, a machine learning algorithm builds a model by examining many examples and attempting to find a model that minimizes loss.



:: How does the model find the "best" parameters ?

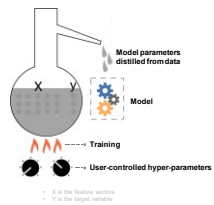
Gradient Descent is one of the most common algorithms to find the good parameters

A Machine Learning model is trained by starting with an initial guess for the parameters (e.g. weights and bias in neural network) and iteratively adjusting those guesses until learning parameters with the lowest possible loss.



:: The model has parameters and hyper-parameters

Iteratively tuning the Hyperparameter so that the model can learn the "best" Parameters from data
The hyper-parameters are specified by the developer/data scientist while parameters are computed from the data via the algorithms.



- Model's parameters are the variables that your chosen machine learning technique uses to adjust to your data. They are internal to the model. They are estimated or learned from data. They are often not set manually by the practitioner.
 - Hyperparameters control how a machine learning algorithm fits the model to the data. Hyper-parameters are specified by the programmer, not computed from the training data, and are often used to tune a model to improve accuracy for a particular data set.
- The examples of hyper-parameters:
- Number of layers, learning-rate in Neural network
 - Number of trees in Random forest

The detailed explanation of the difference will be elaborated later.

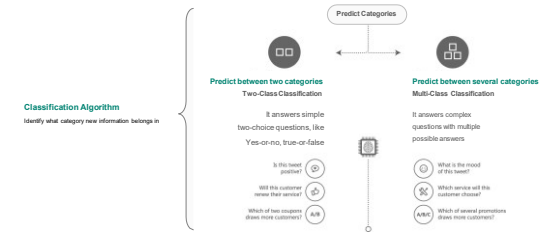
:: Let's take a look at another example for classification problem

Classification algorithm is another major types of machine learning algorithms
Just now we quickly go over the regression model example, now let's move to the example of classification.



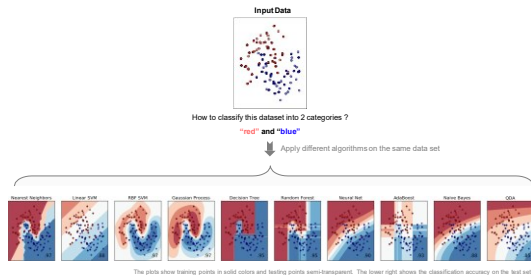
:: Let's take a look at an example for classification problem

Classification algorithm is another major types of machine learning algorithms



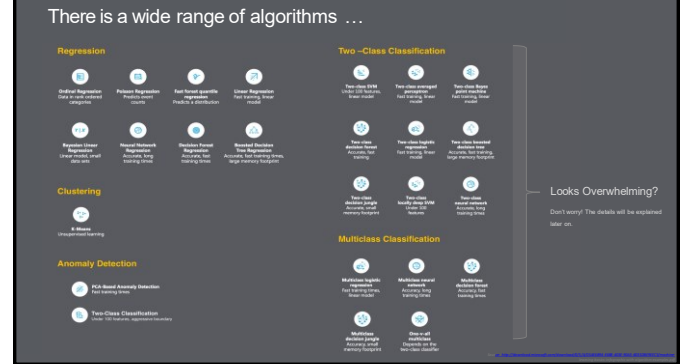
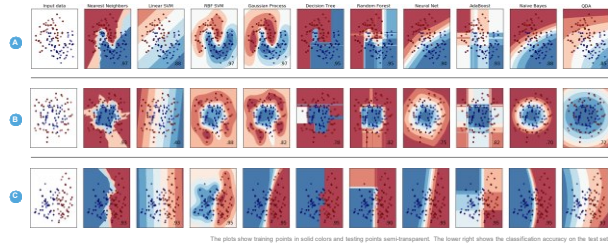
:: Example: how to classify the data points ?

How to classify this dataset into 2 categories ?

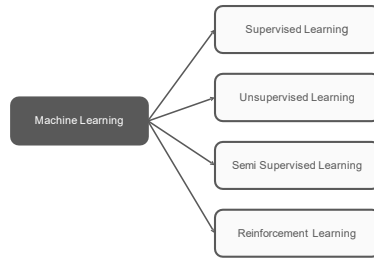


:: Example: how to classify the data points ?

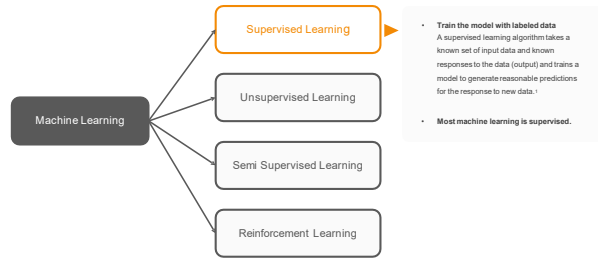
Use different algorithms to classify the data set.
In this example, there are 3 different data sets- A, B, C. You can see how they're classified with different algorithms.



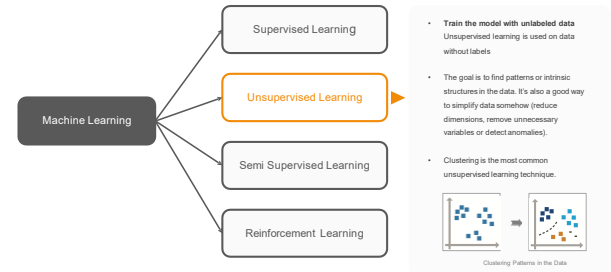
:: The categories of machine learning algorithms



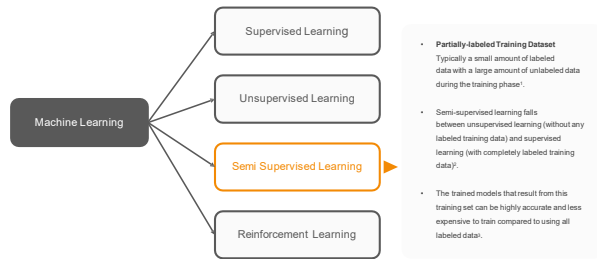
:: The categories of machine learning algorithms



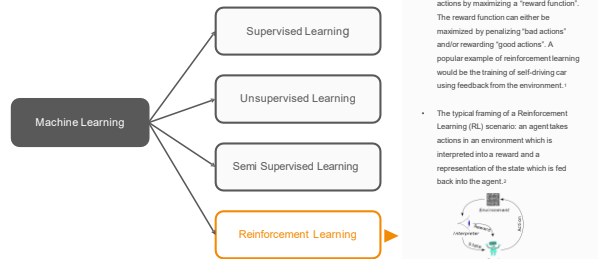
:: The categories of machine learning algorithms



:: The categories of machine learning algorithms

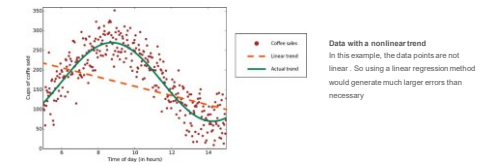


:: The categories of machine learning algorithms



:: You should choose the suitable algorithm

Unsuitable algorithm will result in low accuracy



:: Some of Machine Learning Algorithms

So many algorithms. Which one should I use?

:: Which algorithm should I use ?

Selecting a machine learning algorithm is a process of trial and error
Choosing the right algorithm can seem overwhelming—there are dozens of supervised and unsupervised machine learning algorithms, and each takes a different approach to learning. There is no best method or one size fits all.
Finding the right algorithm is partly just trial and error—even highly experienced data scientists can't tell whether an algorithm will work without trying it out. *

How to select the right algorithm ?
The answer to the question varies depending on many factors, including: *

- The size, quality, and nature of data.
- The available computational time.
- The urgency of the task.
- What you want to do with the data.

Machine Learning

Regression

Dimensionality Reduction

:: Train the model

The modelling is an iterative process of building model which is a combination of data structure, algorithm, and mathematics to captures the relationship between features and target within the dataset.

6 Model Evaluation

:: What's a good model ?

- **Accurate**
Are we making good predictions ?
- **Interpretable**
How easy is it to explain how the predictions are made ?
- **Fast**
How long does it take to build a model and how long does the model take to make predictions?
- **Scalable**
How much longer do we have to wait if we build/predict using a lot more data ?

:: Does your model make good predictions in real situations?

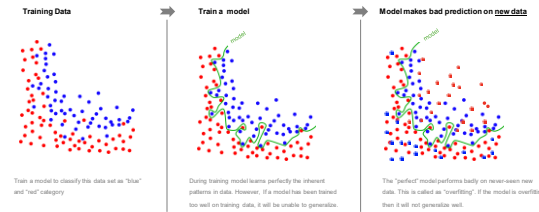
Training Phase

Applying Phase

The model making perfect fitting during training phase might perform very poorly in practice.
The perfect approximation in training phase might be a overfitting problem.

:: The model might not generalize well to unseen new data

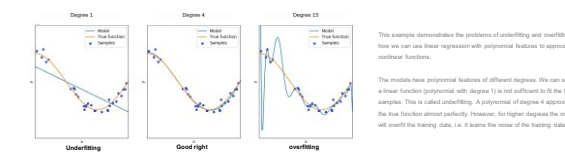
A model's ability to generalize to new data is crucial for the success of a model
Generalization refers to a model's ability to perform well on new unseen data rather than only the training data. If the model is trained too well, it can fit perfectly the random fluctuations or noise in the training data but it will fail to predict accurately on new data.



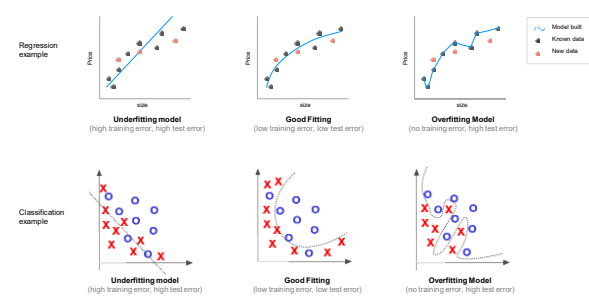
:: About the model fitting

Underfitting vs. Overfitting

- When a model fails to capture important distinctions and patterns in the data, so it performs poorly even in training data, that is called **underfitting**. Underfitting is often a result of an excessively simple model.
- The **overfitting** model performs perfectly on training data but fail to predict well on new data. In other words, it has low/no training error but has high test error.
- Both overfitting and underfitting lead to poor predictions on new data sets.



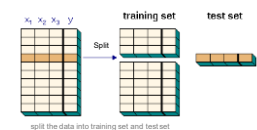
:: About the model fitting



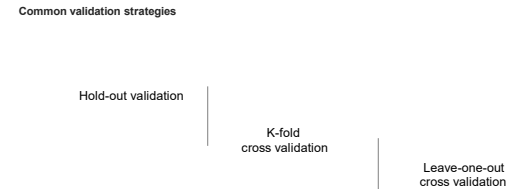
:: Model validation checks how well a model generalizes to new data

Reserve or split a portion of data to validate the model's performance
The fundamental goal of ML is to generalize beyond the data instances used to train models. We want to evaluate the model to estimate the quality of its pattern generalization for data the model has not been trained on.

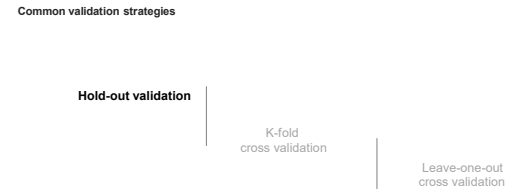
The only way to know how well a model will generalize to new cases is to actually try it out on new cases. It is common practice when performing a (supervised) machine learning experiment to hold out part of the available data as a **test set**. You train your model using the training set, and you test it using the test set. The error rate on new cases is called the generalization error, and by evaluating your model on the test set, you get an estimation of this error. This value tells you how well your model will perform on instances it has never seen before.



:: Model validation strategy



:: Model validation strategy



:: Model validation strategy : Hold-out validation

Split the data into 2 parts - a training set , test set
The Hold-out validation is a type of simple validation.

Usually 70-80% 20-30%

The ML system evaluates predictive performance by comparing predictions on the evaluation data set/test data set with true values (known as ground truth) using a variety of metrics.

However, its evaluation can have a high variance. The evaluation may depend heavily on which data points end up in the training set and which end up in the test set, and thus the evaluation may be significantly different depending on how the division is made!

If you are seeing surprisingly good results on your evaluation metrics, it might be a sign that you are accidentally training on the test set. For example, high accuracy might indicate that test data has leaked into the training set!

:: Model validation strategy : Hold-out validation

Using only two partitions may be insufficient when doing many rounds of hyperparameter tuning
The previous slide introduced partitioning a data set into a training set and a test set. This partitioning enabled you to train on one set of examples and then to test the model against a different set of examples. With two partitions, the workflow could look as follows:

"Tweak model" means adjusting anything about the model you can dream up—from changing the learning rate, to adding or removing features, to designing a completely new model from scratch. At the end of this workflow, you pick the model that does best on the test set.

However, this 2-partitions approach is not recommended although it's widely introduced in many books.

This approach will not be a fair estimate of how well my hypothesis generalizes
Because we fit the parameters using the test set and it gave us the best possible performance on the test set, that's likely to be an overly optimistic estimate of generalization error. So it's no longer fair to evaluate my hypothesis on this test set.

:: Model validation strategy : Hold-out validation - continued

Better approach: Split the data into a training, cross validation and test set
Instead of 2 partitions of data, split data set into 3 parts. Use your cross validation data to select the model and evaluate it on the test data. This approach is recommended by Andrew Ng.

For relatively small data set:
Training data: Usually 60%
Cross Validation data: 20%
Test data: 20%

For big data set:
Training data: 98%
Cross Validation data: 1%
Test data: 1%

This traditional ratio is best practice for small data set. For example, it's good for the smaller data set. For example, 1000, or 10,000 samples.

The traditional splitting ratio (80%-20%-20%) is no longer recommended in big data era. For example, you might have a million samples in total. Then, you can use 98% of the data as training data, and 1% as cross validation data, 1% as test data.

Video: Model Selection And Train Validation Test Sets (13min)
<https://www.youtube.com/watch?v=8tAC0G0D8>

:: Model validation strategy

Common validation strategies

Hold-out validation

K-fold cross validation

Leave-one-out cross validation

:: Model validation strategy : k-fold cross validation

K-fold cross validation: the data set is divided into k equal size subsets.
The data set is divided into k subsets, and the Hold-out validation is repeated k times. Each time, one of the k subsets is used as the test set and the other k-1 subsets are put together to form a training set. Then the average error across all k trials is computed.

Prediction error
= Average Test Error

The advantage of this method is that it matters less how the data gets divided. Every data point gets to be in a test set exactly once, and gets to be in a training set k-1 times. The variance of the resulting estimate is reduced as k is increased. The disadvantage of this method is that the training algorithm has to be run from scratch k times, which means it takes k times as much computation to make an evaluation.

:: Model validation strategy

Common validation strategies

Hold-out validation

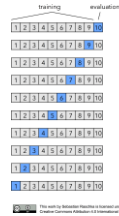
K-fold cross validation

Leave-one-out cross validation

:: Model validation strategy : Leave-one-out cross validation (LOOCV)

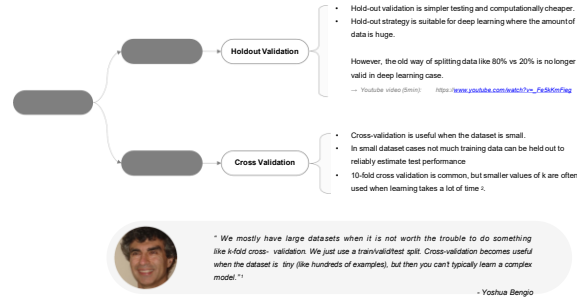
Leave-one-out cross validation: specific case of k-fold cross validation, where $k = n$

Leave-one-out cross validation is K-fold cross validation taken to its logical extreme, with K equal to N, the number of data points in the set. Thus, the learning algorithm is applied once for each instance, using all other instances as a training set and using the selected instance as a single-item test set.

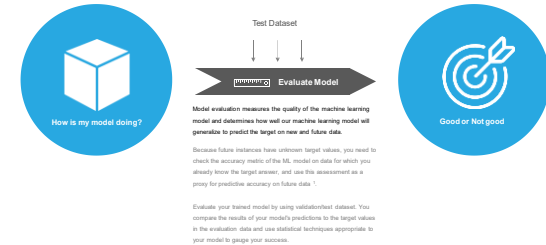


- use one observation as the validation set and the remaining observations as the training set
- This is repeated such that each observation in the sample is used once as the validation data.
- LOOCV is computationally intensive!

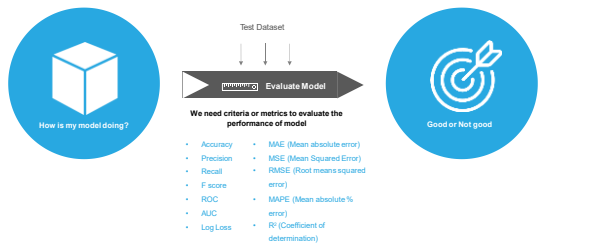
:: Which validation strategy should we use ?



:: Evaluate the model's performance



:: Evaluate the model's performance



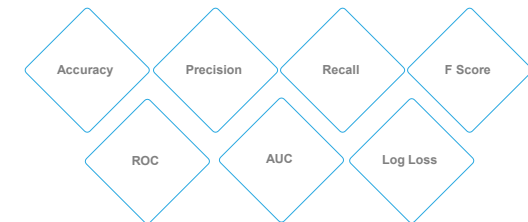
:: Model Evaluation Metrics

Different machine learning tasks have different performance metrics

Classification Model	Regression Model
<ul style="list-style-type: none"> • Accuracy • Precision • Recall • F score • ROC • AUC • Log Loss 	<ul style="list-style-type: none"> • MAE (Mean absolute error) • MSE (Mean Squared Error) • RMSE (Root means squared error) • MAPE (Mean absolute % error) • R² (Coefficient of determination)

evaluation metrics explain the performance of a model

:: Metrics for evaluating classification model



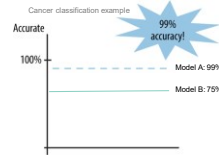
:: What's the accuracy?

Accuracy measures the ratio of correct predictions to the total number of cases evaluated.



:: Accuracy is not a good metric for Skewed Datasets

A predictive model may have high accuracy, but be useless !



Suppose the positive class is only a tiny portion of the observed data. For example, only 1% of patients has true cancer while other 99% of patients don't have any cancers. This is a typical case about data skew.

For such **imbalanced classes**, suppose you build two models.

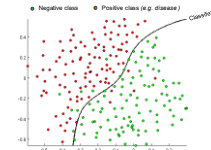
- Model A would achieve 99% accuracy by always predicting "no cancer".
- Model B has 75% accuracy. It correctly predict some cancer classes but also make some misclassification.

Is Model A better than Model B ?

NO ! Model A always classifies incoming data as negative classes (i.e. no cancer) can achieve 99% accuracy but can't identify any cancer cases, therefore it is useless.

Using accuracy is only good for symmetric data sets where the class distribution is 50/50 and the cost of false positives and false negatives are roughly the same*.

:: Confusion matrix visually shows more details of the error

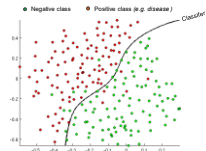


Binary classification example
In this example, the model correctly classifies many classes but also makes some mistakes:

- some positive classes are wrongly classified as negative values.
- some negative classes are wrongly classified as positive classes

How to summarize the prediction result ?

:: Confusion matrix visually shows more details of the error



Binary classification example
In this example, the model correctly classifies many classes but also makes some mistakes:

- some positive classes are wrongly classified as negative values.
- some negative classes are wrongly classified as positive classes

Confusion matrix

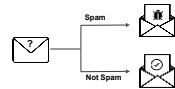
Actual Class	Predicted Class	
	Predicted Value: Positive (+)	Predicted Value: Negative (-)
Actual Value: Positive (+)	TP True Positive	FN False Negative
Actual Value: Negative (-)	FP False Positive	TN True Negative

Confusion Matrix

It is a technique to summarize the performance of classification model. It displays the number of correct and incorrect predictions made by the model compared with the actual classifications in the test data. Main values of this matrix:

- True Positive** - we predicted "+" and the true class is "+"
- True Negative** - we predicted "-" and the true class is "-"
- False Positive** - we predicted "+" and the true class is "-" (Type I error)
- False Negative** - we predicted "-" and the true class is "+" (Type II error)

:: Case study: spam detection



Junk email detection is a binary classification problem. A model is trained to classify the email as "spam" or "not-spam".

Confusion matrix is used to compare the value predicted by the model with the actual value. In this example, there are :

- 12,887 TP (correct positive prediction)
- 98,358 TN (correct negative prediction)
- 358 FP (They are false positive. We predicted them as "spam", but they are not spam)
- 567 FN (we predict them as negative value (i.e. not spam), but they actually are spam)

Actual Class	Predicted Class	
	Predicted Value: Spam (+)	Predicted Value: Non-spam (-)
Actual Value: Spam (+)	TP 12,887	FN 567
Actual Value: Non-spam (-)	FP 358	TN 98,358

Count the number of instances

correct classification (Type I error) wrong classification (Type II error)

:: Example for Accuracy, Precision, Recall

Accuracy

$$\text{Accuracy} = \frac{TP+TN}{TP+FP+TN} = \frac{\text{Correct predictions}}{\text{Total data points}}$$

Precision

$$\text{Precision} = \frac{TP}{TP+FP} = \frac{\text{Correctly Predicted Positive}}{\text{All Predicted Positive}}$$

Recall

$$\text{Recall} = \frac{TP}{TP+FN} = \frac{\text{Correctly Predicted Positive}}{\text{All Real Positive}}$$

Example

Actual Class	Predicted Class	
	Predicted as Spam (+)	Predicted as Non-spam (-)
Actual = Spam (+)	True Positive 12,887	False Negative 567
Actual = Non-spam (-)	False Positive 358	True Negative 98,358

• **Accuracy** = $(TP+TN) / (TP+FP+TN)$

Accuracy = $(2444+18538) / (2444+358+18538) = 77\%$

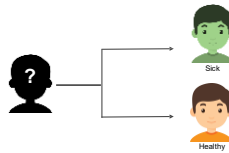
• **Precision** = $\text{True Positives} / (\text{True Positives} + \text{False Positives})$

Precision = $2444 / (2444+358) = 88\%$

• **Recall** = $\text{True Positives} / (\text{True Positives} + \text{False Negatives})$

Recall = $2444 / (2444+567) = 81\%$

:: Case study: Medical diagnosis



Example

	Predicted Class	
	Diagnosed as Sick (+)	Diagnosed as Healthy (-)
Actual = Sick (+)	True Positive 957	False Negative 282
Actual = Healthy (-)	False Positive 1,358	True Negative 7,268

• **Accuracy** = $(TP + TN) / (TP + FP + FN)$
Accuracy = $(957 + 7268) / (957 + 1358 + 7268) = 75.5\%$

• **Precision** = $\text{True Positives} / (\text{True Positives} + \text{False Positives})$
Precision = $957 / (957 + 1358) = 33\%$

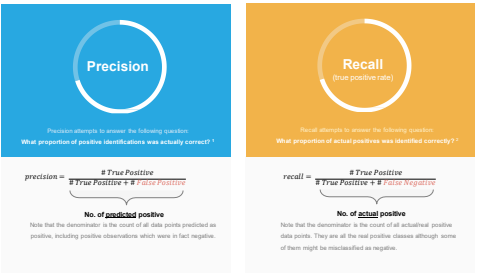
• **Recall** = $\text{True Positives} / (\text{True Positives} + \text{False Negatives})$
Recall = $957 / (957 + 282) = 77\%$

A model is trained to diagnose the people as "sick" or "healthy".

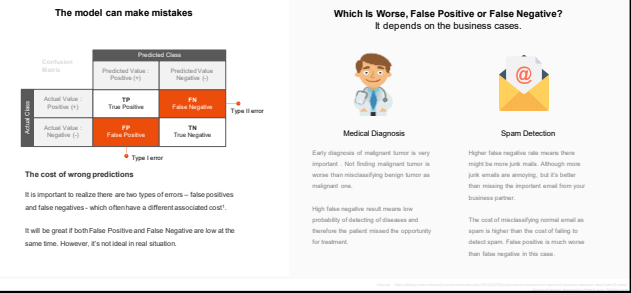
Confusion matrix is used to compare the value predicted by the model with the actual value. In this example, there are:

- 12,987 TP (correct positive prediction)
- 96,356 TN (correct negative prediction)
- 358 FP (False Positive)
- 567 FN (False Negative)

:: Precision vs Recall



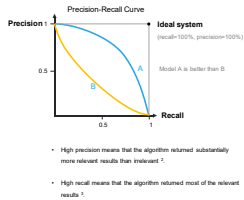
:: About the cost of prediction error



:: Trade-off between precision and recall

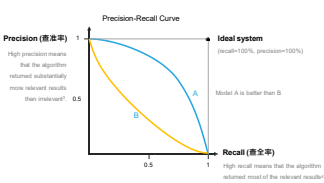
Often, there is an inverse relationship between precision and recall. Increasing precision reduces recall, and vice versa. This is called the precision/recall tradeoff.¹

- Within any one model, you can decide to emphasize either precision or recall. Which one is more important depends on the business needs. In some contexts you mostly care about precision, and in other contexts you really care about recall.
- You can influence precision and recall by changing the threshold of the model. When you select a lower threshold value/cut-off value, then the recall will increase. On the other hand the false positive fraction will also increase, and therefore the true negative fraction/specificity will decrease



:: Trade-off between precision and recall

Change the threshold of a classifier



Let's say we trained a Logistic Regression classifier

- we predict 1 if $\text{f}(x) \geq 0.5$
- we predict 0 if $\text{f}(x) < 0.5$

Case 1: when precision is more important
Suppose we want to predict $y = 1$ (i.e. people have cancer) only if we're very confident. We may change the threshold to 0.7

- we predict 1 if $\text{f}(x) \geq 0.7$
- we predict 0 if $\text{f}(x) < 0.7$

That leads to

- higher precision in this case (if for who we predicted $y = 1$ we are more likely to actually have 1)
- But lower recall (we'll miss more patients that actually have cancer, but we failed to spot them)

Case 2: when recall is more important
Suppose we want to avoid missing too many cases of $y = 1$ (i.e. we want to avoid false negatives). We may change the threshold to 0.3

- we predict 1 if $\text{f}(x) \geq 0.3$
- we predict 0 if $\text{f}(x) < 0.3$

That leads to

- Higher recall (we'll correctly flag higher fraction of patients with cancer)
- Lower precision (and higher fraction will turn out to actually have no cancer)

:: Metric: F score

How to compare precision/recall numbers and decide which algorithm is better ?

	Precision(P)	Recall(R)	Which is best?
Algorithm 1	0.5	0.4	1
Algorithm 2	0.7	0.1	2
Algorithm 3	0.02	1	3

- By varying threshold parameter we will get different precision and recall
- Improving recall will lead to worse precision
- Improving precision will lead to worse recall
- How to pick the threshold/cut-off value?
- Is there a way to automatically choose the threshold for us?



Which algorithm is better?

- Now have two numbers and need to choose which one to prefer
- F1 score helps to decide since it's just one number

:: Metric: F score

F score combines precision and recall into one measure

	Precision(P)	Recall(R)	F ₁ Score
Algorithm 1	0.5	0.4	0.444 ✓
Algorithm 2	0.7	0.1	0.175
Algorithm 3	0.02	1	0.0392

• The F₁ score is a good way to summarize the evaluation of performance in a single number. It combines precision and recall into one measure.

• Bigger F₁ score is better. The ideal F₁ score value is 1. The worst value is 0. In this example, Algorithm 1 has highest F₁ score, so it's the best one.

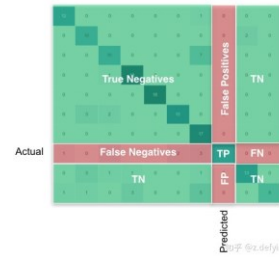
A combined measure: F score

$$F_1 = 2 \frac{PR}{P+R} = 2 \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

• It is often convenient to combine precision and recall into a single metric called the F1 score, in particular if you need a simple way to compare two classifiers.

• The F₁ Score is the weighted average of Precision and Recall. Therefore, this score takes both false positives and false negatives into account.

:: Metric: Multi-class Confusion Matrix



- TP: the sum of diagonal value
- FP for specific class: the sum of corresponding column minus the value of corresponding tp
- FN for specific class: the sum of corresponding row minus the value of corresponding tp
- TN for specific class: the sum of the matrix minus the sum of (tp + fp + fn) of the corresponding class

:: Metric: Confusion Matrix exercise (10 min)

信用卡欺诈行为判断：一共有500条数据，其中正常行为400条，欺诈行为100条
 模型A：预测正常行为为270条，整体准确率66%
 模型B：预测正常行为为400条，整体准确率80%

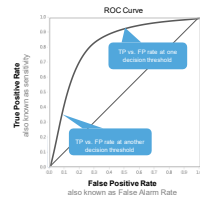
Q: 请问应该采取哪个模型的结果？



:: Metric: ROC Curve

ROC curve plots **True Positive Rate** vs. **False Positive Rate** at different classification thresholds

ROC Curves shows the tradeoff between true positive rate and false-positive rates of classification algorithms. In other words, ROC shows you how many correct positive classifications can be gained as you allow for more and more false positives.



Based on confusion matrix, we can calculate True Positive Rate and False Positive Rate.

• **True Positive Rate (Recall)**
 It is sometimes referred to as the hit rate - what percent of the actual positives does the classifier get right.

$$TPR = \frac{\# \text{True Positive}}{\# \text{True Positive} + \# \text{False Negative}}$$

• **False Positive Rate**
 It is the ratio of negative instances that are incorrectly classified as positive

$$FPR = \frac{\# \text{False Positive}}{\# \text{False Positive} + \# \text{True Negative}}$$

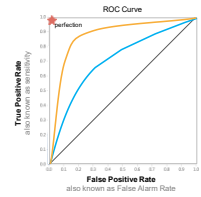
True/False: Watch: How to plot the ROC curve ? (video)

https://youtu.be/47T03L4G7Y0

:: Metric: ROC Curve

The closer this curve is to the upper left corner, the better the classifier's performance is

Curves that are close to the diagonal of the plot, result from classifiers that tend to make predictions that are close to random guessing¹. The perfect classifier that makes no mistakes would hit a true positive rate of 100% immediately, without incurring any false positives—this almost never happens in practice².



In this example

- Classifier A is better than B
- Classifier B is better than Random Guessing

Main advantages of ROC:

- ROC Curves are insensitive to class distribution/unbalanced datasets.
- If the proportion of positive to negative instances changes, the ROC Curve will not change.

Note that no classifier should be in the lower right triangle of a ROC graph.

This represents performance that is worse than random guessing!

1. Source: https://www.kdnuggets.com/2011/05/roc-curves.html

2. Source: https://www.kdnuggets.com/2011/05/roc-curves.html

:: Metric: AUC (area under the curve)

AUC is (arguably) the best way to summarize model's performance in a single number

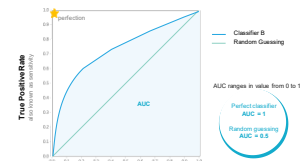
• Compare multiple models by a single number

The ROC curve is not just a single number. It's a whole curve. It provides nuanced details about the behavior of the classifier, but it's hard to quickly compare many ROC curves to each other. In particular, if one wants to employ some kind of automatic hyper-parameter tuning mechanism, the machine would need a quantifiable score instead of a plot that requires visual inspection. The AUC is one way to summarize the ROC curve into a single number, so that it can be compared easily and automatically.

• Higher AUC will be the better

Using the AUC metric you can quickly compare multiple learning models, when comparing two models, the one with a higher AUC will be the better one regardless of the threshold setting. Compared to the statistical measures of accuracy, precision, recall and F1 score, AUC's independence of threshold makes it uniquely qualified for model selection.

On the other hand, unlike accuracy, precision, recall and F1 score, AUC does not tell us what performance to expect from the model for a given threshold setting, nor can it be used to determine the optimal value for threshold. In that regard it doesn't take away the need for the other statistical measures.



False Positive Rate
 also known as False Alarm Rate

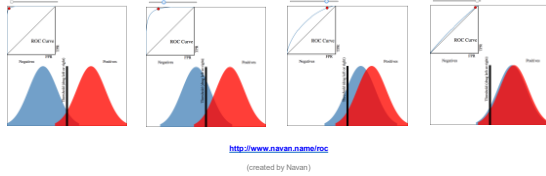
In practice, if you have a "perfect" classifier with an AUC of 1.0, you should be suspicious, as it likely indicates a bug in your model. For example, you may have overfit to your training data, or the label data may be replicated in one of your features.

1. Source: https://www.kdnuggets.com/2011/05/roc-curves.html

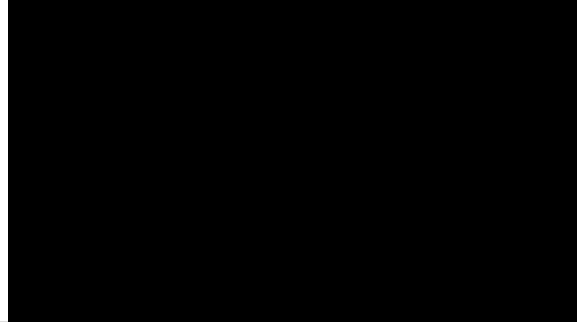
2. Source: https://www.kdnuggets.com/2011/05/roc-curves.html

:: Good demo: ROC curve & AUC

An interactive visualization that may help you better understand ROC curves and AUC



:: Good demo: Acc, Precision, Recall



:: Good demo: ROC curve & AUC



:: Metrics for evaluating regression model

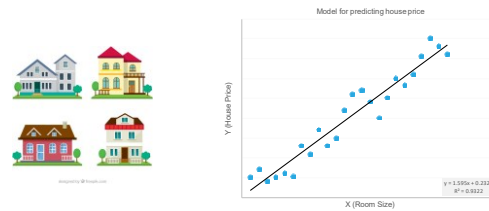


The metrics returned for regression models are generally designed to estimate the amount of error. A model is considered to fit the data well if the difference between observed and predicted values is small. However, looking at the pattern of the residuals (the difference between any one predicted point and its corresponding actual value) can tell you a lot about potential bias in the model.

:: Linear regression errors example

Linear Regression example

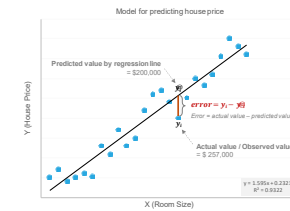
Suppose you want to create a simple model for predicting the house price just based on the room size. A linear regression model can be built for predicting the price.



:: Linear regression errors example

The prediction error

You can compare the observed value with the predicted value. Such error can be used to evaluate the model's performance.



Linear Regression Errors

There are difference between the actual value and predicted value from the model. That's the prediction error.

$$Error = y_i - \hat{y}_i$$

In this example, error for this data point:

$$Error = 200,000 - 257,000 = -57,000$$

To find the "best fit" line, the goal is to minimize the sum of the errors of all data points. However, the negative values can cancel out the positive values if only calculating the sum. So some techniques can be used to address this issue. See next slides for more details.

Common metrics for measuring the regression errors

MAE
Mean Absolute Error

MSE
Mean Squared Error

RMSE
Root Mean Squared Error

MAPE
Mean Absolute Percentage Error

Lower error values mean the model is more accurate in making predictions

Metric for regression model: MAE

MAE is the mean of the absolute value of the errors

Predicting the house price

Y (House Price)

X (Room Size)

$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$

$y = 1.589x + 0.2201$
 $R^2 = 0.9322$

Mean Absolute Error (MAE)

- MAE measures how close the predictions are to the actual outcomes.
- It's the average sum of absolute difference between the actual value the predicted values for all data points. Because of the absolute value, the negative and positive errors don't cancel out each other
- The best regression is the one that minimizes MAE. Thus, a smaller score is better!

Metric for regression model: MSE

MSE is the mean of the squared value of the errors

Predicting the house price

Y (House Price)

X (Room Size)

$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$

$y = 1.589x + 0.2201$
 $R^2 = 0.9322$

Mean Square Error (MSE)

- MSE measures the average sum of squares of the difference between the actual value the predicted values for all data points. Because of the square, the negative values don't cancel positive values and it also amplifies the impact of the errors.
- Because of the square, large errors have relatively greater influence on MSE than do the smaller error.
- The best regression is the one that minimizes MSE. Thus, a smaller score is better!

Metric for regression model: RMSE

RMSE is the square root of the mean of the squared errors

Predicting the house price

Y (House Price)

X (Room Size)

$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$

$y = 1.589x + 0.2201$
 $R^2 = 0.9322$

Root Mean Square Error (RMSE)

- RMSE is the most popular evaluation metrics used in regression problems. It has same unit with "Y".
- RMSE may be the most common metric, but it has some problems. RMSE is highly affected by outlier values.
- RMSE more aggressively punishes big errors than small ones. This means the RMSE should be more useful when large errors are particularly undesirable.

The RMSE is more sensitive to outliers than the MAE. But when outliers are exceptionally rare this is not always correct. The RMSE performs very well and is generally preferred!

- A smaller RMSE value is better!

Metric for regression model: MAPE

The MAPE (Mean Absolute Percent Error) measures the size of the error in percentage terms

Predicting the house price

Y (House Price)

X (Room Size)

$MAPE = \frac{100}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right|$

$y = 1.589x + 0.2201$
 $R^2 = 0.9322$

Mean Absolute Percentage Error (MAPE)

- MAPE functions best when there are no extremes to the data (including zeros).
- In spite of the popularity in measuring the forecasting error, MAPE has major drawbacks in practical application and it should be used cautiously.

With zeros or near-zeros, MAPE can give a distorted picture of error. The error on a near-zero item can be relatively high, creating a distorted view of the model error rate when it is averaged in. For forecasts of items that are near or at zero volume, Symmetric Mean Absolute Percent Error (SMAPE) is a better measure!

- A smaller MAPE value is better

MAPE is one of the most popular measures for forecasting time series error.

Other metrics for evaluating the regression model

R^2

R^2_{adj}

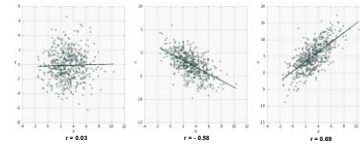
Bigger value is better

:: Let's learn the key concept one by one



:: At first let's take a look at r

r is also called as Pearson correlation coefficient



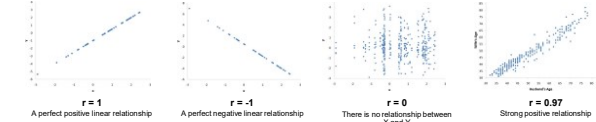
Pearson correlation coefficient, also referred to as the Pearson's r , is a measure of the linear correlation between two variables X and Y .

By examining the coefficient values, you can infer something about the strength of the relationship between the two variables, and whether they are positively correlated or negatively correlated.

- It has a value between +1 and -1.
- 0 is no linear correlation

:: At first let's take a look at r

Pearson's r measures the strength of the linear relationship between two variables.



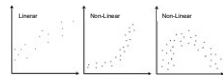
- A positive correlation indicates a relationship between x and y measures such that as values of x increase, values of y also increase.
- A negative correlation indicates the opposite—as values of x increase, values of y decrease

:: At first let's take a look at r

How to calculate Pearson's r ?

Step 1: Determine linearity

The first step in studying the relationship between two continuous variables is to draw a scatter plot of the variables to check for linearity. The correlation coefficient should not be calculated if the relationship is not linear.



Step 2: Clean data

You must remove or fill in missing values, remove or clip outliers, and ensure that the columns have the proper data type.

Step 3: Calculate the coefficient

$$r = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2 \sum (y - \bar{y})^2}} \quad \text{Or} \quad r = r_{xy} = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \sqrt{n \sum y_i^2 - (\sum y_i)^2}}$$

where:
• n is the number of samples
• \bar{x} is the mean of x
• \bar{y} is the mean of y
• \sum is the sum of the values in the column

Example

x	y	x ²	y ²	xy
1	2	1	4	2
2	3	4	9	6
3	4	9	16	12
4	5	16	25	20
5	6	25	36	30
6	7	36	49	42
7	8	49	64	56
8	9	64	81	72
9	10	81	100	90
10	11	100	121	110
11	12	121	144	132
12	13	144	169	156
13	14	169	196	182
14	15	196	225	210
15	16	225	256	240
16	17	256	289	272
17	18	289	324	306
18	19	324	361	342
19	20	361	400	380
20	21	400	441	420

Video: Correlation Coefficient (3 min)
<https://www.youtube.com/watch?v=26QnT78K>

:: R^2 - coefficient of determination

R^2 is a standard way of measuring how well the model fits the data



• The coefficient of determination, denoted R^2 or r^2 and pronounced "R squared"
It is the square of the Pearson correlation coefficient r . So, for example, a Pearson correlation coefficient of 0.8 would result in a coefficient of determination of 0.64, (i.e., $0.8 \times 0.8 = 0.64$).

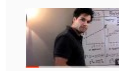
• R^2 describes the proportion of the variance/variation in the dependent variable y explained by the regression model.

For example, $R^2 = 0.93$ or 93% means that approximately 93% of the variation the dependent variable y can be explained by the variation in the variable x .

Still confused? Please watch the video for detailed explanation. →



Video: What does r squared tell us? What does it all mean



Youtube Video (3 min)

<https://www.youtube.com/watch?v=18GdW0C-2>



:: R^2 - coefficient of determination

In general, the higher the R -Squared, the better the model fits the data

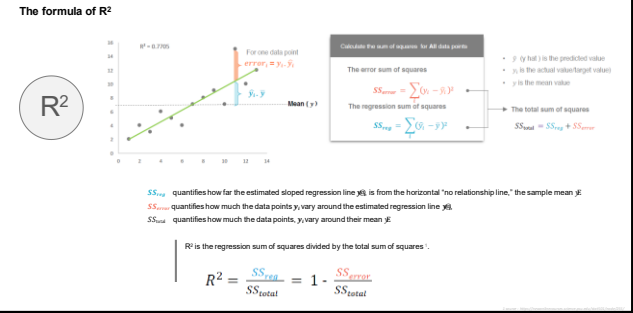


R^2 gives us some information about the goodness of fit of a model. In other words, R^2 measures how well the regression line approximates the real data points.

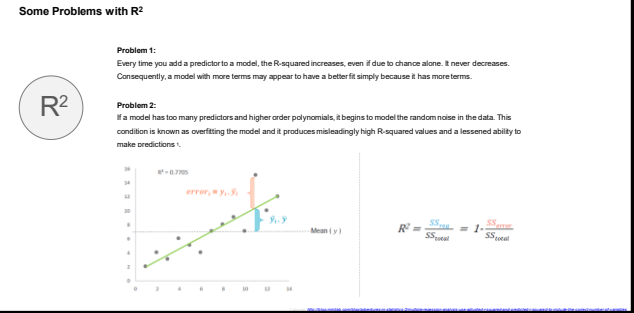
R^2 is a value between 0 and 1.

- 0 means the model is random (explains nothing).
- 1 means there is a perfect fit.
- However, caution should be used in interpreting R^2 values, as low values can be entirely normal and high values can be suspect.

:: R² - coefficient of determination



:: R² - coefficient of determination



:: Adjusted R²

