

A Foray Into Machine Learning

Q & A on the main concepts and terminology

Huascar Sanchez

github.com/hsanchez

Home Research Lab

June 13, 2020

The views expressed do not necessarily reflect the position of my employer.

Q. What is Machine Learning?

*"Machine Learning (or **ML**) is fitting a function to examples¹ and using that function to generalize and make predictions about new examples."*

Derek Jedamski, GitHub

Machine Learning, by large, falls into three categories:

- Supervised learning
- Unsupervised learning
- Reinforcement learning

And solves two types of problems: Classification and Regression.

¹An example, or a sample, is a data point that belongs to some data

Q. What is Supervised Learning?

In **Supervised Learning** (or **SL**), you are given a bunch of examples and their labels (e.g., A or B) and the goal is to classify, when you are given a new example, to which label we should assign the new example.

You could think of these labels as the names of the **classes or clusters** to which certain portions of the data belong.

Q. Examples of supervised learning?

Classification (pattern recognition):

Speech recognition, machine translation, ...

Medical diagnosis: often, variables are missing (tests are costly).

Regression (the labels to be predicted are continuous)

Predict the price of a car from its mileage, ...

Kinematics of a robot arm: predict workspace location from angles

Q. Classification vs Regression?

Classification is the task of predicting a discrete class label.

Regression is the task of predicting a continuous quantity.

There's some overlap between classification and regression algorithms; for example

A classification algorithm may predict a continuous value, but the continuous value is in the form of a probability for a class label.

A regression algorithm may predict a discrete value, but the discrete value is in the form of an integer quantity.

Q. Example of supervised learning algorithm?

Support vector machines or SVM, is a ML algorithm that takes some data as input and returns as output an optimal separating *hyperplane* between pieces of data; e.g., a plane separating A from B.

- In this Ex, the data sit in some high dimensional space, and the idea is to construct a plane that maximizes the margins² between the plane and the data. If a new datum sits closer to one area of the data, say A, then we assign this new datum to A.

(For historical reasons) This algorithm is called **support vector machines** because *the vectors that lie on the margin of the plane* are called the **support vectors**.

²distance between points closest to the line and the actual line

Q. What more can you say about SVM?

In summary, SVM is a method for constructing a device to discriminate. If we're having a supervised learning problem then this method gives me an optimal form of discrimination (i.e., optimal separating hyperplane³).

SVM works on both linearly and non-linearly separable data. In latter case, it use the kernel trick⁴ to convert these data to linearly separable data in a higher dimension.

This transformation is called kernel.

³A hyperplane in an n -dimensional Euclidean space is a flat, $n - 1$ dimensional subset of that space that divides the space into two disconnected parts.

⁴adds one more dimension (called z -axis) – governed by the constraint $z = x^2 + y^2$ and z is the squared distance of the points from origin.

Q. Advantages and Disadvantages SVM?

Advantages

SVM Classifiers offer **good accuracy and perform faster prediction** compared to Naive Bayes algorithm. They also **use less memory** because they use a subset of training points in the decision phase. SVM works well with a clear margin of separation and with high dimensional spaces.

Disadvantages

SVM is **not suitable for large data sets** because of its high training time and it also takes more time in training compared to Naive Bayes. It **works poorly with overlapping classes** and is also **sensitive** to the type of **kernel** used.

Q. SVM hyperparameters?

SVM has a few hyperparameters⁵; however, its most popular include:

- 1 Kernel parameter (data transformation; kernel trick).
- 2 C regularization parameter (misclassification penalty).
- 3 γ (*Gamma*) parameter (spread of kernel/decision boundary).

⁵**Hyperparameters** are the properties that govern the entire ML training process. They are external to the model and whose values cannot be estimated from data. We use them to help estimate model parameters.

Q. SVM's kernel hyperparameter?

Kernel is a transformation that “transforms” a given data input into a desired form; e.g., from non-linearly separable data to linearly separable data in a higher dimension.

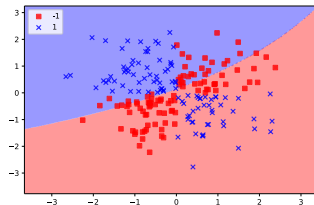
There are various types of kernels such as linear, polynomial, and *radial basis function* (RBF). Polynomial and RBF⁶ are useful for non-linear hyperplanes.

⁶Unlike the polynomial kernel which looks at d extra dimensions, RBF expands into an infinite number of dimensions; enabling the inner product of data points that have an infinite number of dimensions.

Q. SVM's C (regularization) hyperparameter?

C is the penalty for misclassifying data points:

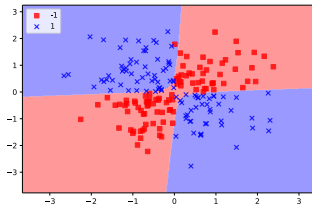
- When C is small, the classifier is okay with misclassifying data points; if it's too small it can lead to underfitting (high bias, low variance).



Q. SVM's C (regularization) hyperparameter?

C is the penalty for misclassifying data points:

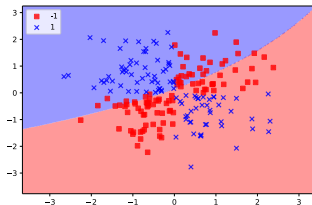
- When C is small, the classifier is okay with misclassifying data points; if it's too small it can lead to underfitting (high bias, low variance).
- When C is large, the classifier is heavily penalized for misclassifying data points and therefore it bends over backwards avoiding any misclassified data points; if it's too large it can lead to overfitting (low bias, high variance).



Q. SVM's γ hyperparameter (of Radial Basis Function (RBF) Kernel)?

γ (Gamma) is the “spread” of the RBF kernel & hence the decision region.

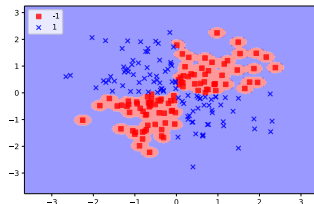
- When γ is low, the “curve” of the decision boundary is very low (like an arch) and thus the decision region is very broad.



Q. SVM's γ hyperparameter (of Radial Basis Function (RBF) Kernel)?

γ (Gamma) is the “spread” of the RBF kernel & hence the decision region.

- When γ is low, the “curve” of the decision boundary is very low (like an arch) and thus the decision region is very broad.
- When γ is high, the “curve” of the decision boundary is high; meaning the decision boundary mostly depends on individual data points, creating islands of decision-boundaries around data points (may overfit).



Q. What is Unsupervised Learning?

In **Unsupervised Learning** (or **UL**), you are given a bunch of data and you are not told they fall naturally into clusters, and also you are not told what these clusters are.

The goal is identify the clusters within data, how many clusters there are, and then be able to assign new things to these different clusters.

Q. Examples of unsupervised learning?

Learning associations

- Basket analysis. Let $\Pr(Y | X)$ the probability that a customer who buys X also buys Y – estimated from past purchases. If $\Pr(Y | X)$ is 0.7 then associate X to Y ($X \rightarrow Y$); meaning when someone buys X , recommend Y .

Clustering (group similar data points).

Density estimation (where are data points likely to lie?)

Feature selection (keep only useful features).

Outlier/novelty detection.

Q. Can you give an example of Unsupervised Learning?

Principal Component Analysis (or PCA). In general terms, PCA is about finding the underlying patterns of the data (i.e., those simply expressible correlations between your data⁷) and also giving you a method for data compression.

For instance, if your data points in your data are highly correlated with each other and there are only a few forms of correlations, then your covariance matrix is low rank and can be approximated⁸ by only a few principal components.

PCA is all about diagonalizing this covariance matrix.

PCA is an exercise in linear algebra on very high dimensional vector spaces.

⁷Represented as a covariance matrix

⁸All your vectors can be written as the sum of a few w s.

Q. Unpacking PCA?

Principal Component Analysis (or **PCA**) is a classical UL algorithm.

In **PCA**, the way this works, we construct a *covariance matrix*, and this covariance matrix is just the following object: $C = \sum_j \vec{v}_j \vec{v}_j^+$, where \vec{v}_j^T is the transpose of \vec{v}_j .

- In other words, we construct C from the data by taking these vectors \vec{v}_j and multiply them by their transpose \vec{v}_j^+ .

Then, we diagonalize C and say $C = \sum_k P_k \vec{\omega}_k \vec{\omega}_k^{+9}$, P_k is piece of the data with size k , and $\vec{\omega}_k$ are the set of vectors you need to find.

And if only a small number of $P_k \gg 0$, then C is effectively low-rank, and the corresponding $\vec{\omega}_k$ are the principal components. In other words, you need find the eigenvectors that have the largest eigenvalues – i.e., your principal components.

⁹ C can be decomposed into a product of matrices involving eigenvalues and eigenvectors

Q. When to use PCA?

When to use it

PCA should be used if one to figure out if there are latent features driving the patterns in the data. E.g., The big shots at SRI International.

Dimensionality reduction (and feature selection). E.g., helps you visualize high dimensional data, helps you reduce noise in your data, make other ML algos work better (regression, classification) because fewer inputs.

When not to use it

PCA is not suitable in many cases: For example, if all the components of PCA have quite a high variance, there is no *good* universal stopping rule that allows you to discard some exact k Principal Components, meaning no good data compression.

Not good when working with fine-grained classes¹⁰.

¹⁰hard-to-distinguish object classes

Q. What is Reinforcement Learning?

Reinforcement learning is one of three basic machine learning paradigms, alongside supervised learning and unsupervised learning.

In reinforcement learning an algorithm learns to do something by being rewarded for successful behavior and/or being punished for unsuccessful behavior. No supervised output but delayed reward. **Ex:** playing chess or a computer game, robot in a maze.

Q. Can you give an example of Reinforcement Learning?

Policy Gradient (PG)

In this method, we have the policy π that has a parameter θ . This π outputs a probability distribution of actions.

Then we must find the best parameters (θ) to maximize (optimize) a score function $J(\theta)$, given the discount factor γ and the reward r .

Main steps:

- Measure the quality of a policy with the policy score function

- Use policy gradient ascent to find the best param that improves the policy.

Data representation

Q. How do you represent data in ML?

In general, the given data is expressed in a form of a bunch of vectors $\vec{v}_j \in \mathbb{R}^d$ that belong to some high dimensional vector space.

For instance, in image recognition, the vector¹¹ of an each image is a set of pixels¹² (i.e., a pixelated version of the image).

If you have a notion of distance $\Delta(\vec{v}_i, \vec{v}_j)$, then you can compare which vectors are close to each other in this high dimensional vector space; e.g., the norm $\|\vec{v}_i - \vec{v}_j\|^2$.

¹¹a.k.a., feature vector; it could be numerical (e.g., height of tree) or descriptive (e.g., eye color)

¹²each entry in the vector (e.g., pixel) represents a feature.

Q. What is it important about measuring distances?

Because the shorter the distance between two feature vectors the closer in character are the two samples they represent.

There are many ways to measure distance:

- **Manhattan distance** (L^1 norm). The sum of the absolute values of the different between entries in the vector. (Preferred dist. when dealing with high dimensional data.)
- **Euclidean distance** (L^2 norm). Square the distances between vector entries, sum these and square root.
- **Cosine similarity**. Cosine of the angle between two vectors¹³. Just take the *dot* product of two vectors and divide by the two lengths.

¹³Two vectors might be similar if they are pointing in the same direction even if they are of different lengths.

Q. Any other issues we should care to learn?

The curse of dimensionality. This phenomena involves data in high dimensions.

Ex. Suppose you are working in M dimensions, that is you data has M features. And suppose you have N data points. Having a large number of data points is good, the more the merrier. BUT what about number of features?

Think how these N data points might be distributed over M dimensions. Suppose that the numerical data for each feature is 0 or 1. Therefore, there will be 2^M possible combinations.

If N is less than 2^M then you run the risk of having every data point being on a different corner of the M -dimensional hypercube.

Building a Machine learning model

Q. What does this process looks like?

- 1 **Explore & clean the data** to learn about feature relationships, the shape of the data, correlations to the target variable, missing values, and then to make the signal clear and easy for a model to pick up.
- 2 **Split data into train/validation/test data sets** to ensure model does not memorize the examples and learns the underlying pattern.
- 3 **Fit a basic model & evaluate** using cross validation (e.g., 5-fold) on the training set to set a baseline.
- 4 **Tune hyper parameters using GridSearchCV** to explore different hyper parameters combinations and select those hyper parameter settings that generate the best model¹⁴.
- 5 **Evaluate best models on validation set** to select the best model after refitting the top model of each algorithm.
- 6 **Select and evaluate the best model on the test set** to get an unbiased view of the model performance on completely unseen data.

¹⁴those models that beat our baseline

Measuring success
(**hint:** we are in step 2 of process for
building ML models)

Q. What does it mean to measure success of the model?

Q. how do we make sure the model is learning the underlying pattern and not just memorizing the examples?

→ **A.** we split our data into 3 data sets: **training, validation, and testing.**

(Again), this step is about making sure¹⁵ the model **learns** the underlying pattern (or correlation) and be able to **generalize** and **make predictions** about future examples.

¹⁵We don't know how well the model will generalize for we don't have any additional data to test this.

Q. Training, validation, and test data sets?

We typically split a data set into three separate data sets; training (60%), validation (20%), and testing (20%) data sets.

```
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(features, labels,
                                                    test_size=0.4, random_state=42)

X_test, X_val, y_test, y_val = train_test_split(X_test, y_test,
                                                test_size=0.5, random_state=42)
```

Training data, or examples from which the model will learn the general patterns.

Validation data, or examples we use to select the best model (algorithm and hyper-parameter settings).

Test data, or examples we use to provide an unbiased evaluation of what the model will look like in its real environment.

Q. How do we use the training, validation and test sets?

1. You start train your ML algorithm on the training data, evaluate it using the validation data, then at this point:

IF none of your models are any good based on the performance on the validation set, then you need to revisit the training phase and consider some new variables or new models. (Go back to step 1.)

ELSE (performance is quite good) select your best model and pass it onto the testing phase. (Go to step 2.)

2. You then evaluate the best model on the test set.

IF the performance is what you expected, that model is ready to go. ELSE go back to the drawing board. (Got to step 1.)

Q. Supporting this process with cross-validation? Holdout test set?

K-Fold Cross-Validation: Data is divided into k subsets and the **holdout method** is repeated k times. Each time, one of the k subsets is used as the test set and the other $k - 1$ subsets are combined to be used to train the model.

Holdout method uses the **Holdout test set**, a generalization of the test set. The holdout test set is just any data set that was not used in fitting a model (data set aside for evaluating the model's ability to generalize).

Q. 5-Fold Cross-Validation Example?

E.g., 5-fold CV: Given 10,000 examples, you partition into 5 buckets, each consisting of 2,000 examples¹⁶.

On trial 1, we set aside the last bucket (holdout test set) and the other 4 buckets are the training set, then compute its predictive performance. **On trial 2**, we set aside the penultimate bucket, and the other 4 buckets are the training set. This will be similar for the other 3 **trials**.

At the end we *average* the performance of model over the many trials.

¹⁶this partitioning is done thru sampling without replacement; i.e., no single example will appear in two different subsets and all original 10,000 are still accounted for in these subsets

Q. Putting it all together (An evaluation framework)

Our cohesive framework for evaluating our model, by large, consists of two components:

1. **Evaluation metrics:** how are gauging the accuracy of the model?
2. **Process (how to split the data):** how do we leverage a given data set to mitigate the likelihood of overfitting and underfitting.

Evaluation metrics:
There isn't a one-size-fits-all metric.

Q. What are the evaluation metrics that we'll use?

The metric(s) chosen to evaluate a ML model depends on various factors:

Is it a regression or a classification task? E.g., MSE vs Accuracy.

What is the business objective? E.g., precision vs recall.

What is the distribution of the target variable?

Examples of other metrics: Mean Absolute Error (MAE), Mean Squared Error (MSE), R-squared, Confusion Matrix and related metrics (Precision, Recall, Accuracy).

Q. Measuring classifier performance?

Binary classification:

Accuracy¹⁷ in %: $\#(\text{classified correctly}) / \#(\text{total examples})$

Precision: $\#(\text{classified correctly as A}) / \#(\text{total examples predicted in A})$

Recall: $\#(\text{classified correctly as A}) / \#(\text{total examples known in A})$. We can always achieve perfect recall by returning the entire dataset (which will contain many irrelevant examples)

¹⁷or classification error: $\#(\text{misclassified}) / \#(\text{total examples})$ (1 - Accuracy)

Q. Measuring classifier performance?

Binary classification:

Receiver operating curve (ROC): the pair values (false positive rate or FPR, true positive rate or TPR) as a function of a threshold $\theta \in [0, 1]$. An ideal classifier is at $(0, 1)$ (top left corner). Diagonal ($FPR = TPR$): random classifier. *This is the worst we can do.* Any classifier that is below the diagonal can be improved by flipping its decision.

Area under the curve (AUC): It reduces the ROC to a number. Ideal classifier: $AUC = 1$. ROC and AUC allow us to compare classifiers over different loss conditions, and choose a value of θ accordingly. Often, there is not a dominant classifier.

Q. Measuring classifier performance?

$K > 2$ classes:

Again, the most basic measure is the classification error.

Confusion matrix: $K \times K$ matrix where entry (i, j) contains the number of instances of class C_i that are classified as C_j .

It allows us to identify which types of misclassification errors tend to occur, e.g. if there are two classes that are frequently confused. **Ideal classifier:** the confusion matrix is diagonal.

Q. What is a confusion matrix?

A **Confusion Matrix** is a simple way of understanding how well an algorithm is doing at classifying data. It is just the idea of **false positives** and **false negatives**

True class	Predicted class		
	Positive	Negative	Total
Positive	tp: true positive	fn: false negative	p
Negative	fp: false positive	tn: true negative	n
Total	p'	n'	N

Q. Can you explain what precision and recall are?

Recall is a measure of completeness or quantity, whereas

Precision is a measure of exactness or quality:

$$\text{Recall} = \frac{TP}{TP + FN}$$

What proportion of actual positives
was identified correctly?

$$\text{Precision} = \frac{TP}{TP + FP}$$

What proportion of positive identifications
was actually correct?

High precision means your algorithm has returned substantially *more relevant results than irrelevant ones*, while **high recall** means it has returned *most of the relevant results*.

Q. Can you explain what are false positives and false negatives?

False positives and **false negatives**, technically referred to as type I error and type II error respectively.

False positives are incorrect classifications of the presence of a condition when it is actually absent. A false positive is when you reject a true null hypothesis.

False negatives are incorrect classifications of the absence of a condition when it is actually present. A false negative is when you accept a false null hypothesis.

Q. Provide examples when false positives are more important than false negatives, false negatives are more important than false positives and when these two types of errors are equally important

Case 1, **Airport security**. Ensuring that truly dangerous items like weapons cannot be brought on board an aircraft. Getting false positives is better than getting false negatives: missing cases of actual weapons could lead to dangerous situations.

Case 2, **Cancer screening**. Even though false-positive results could create anxiety and lead to unnecessary and invasive follow-up tests like biopsies, missing cases of actual cancer could lead to delays in treatment that negatively affect somebody's life.

Case 3, **General forecasting**. Measuring success of testing cases of covid-19 in the US. Tracking rate of false positives and false negatives seems to make sense.

Q. Can you list other metrics derived from the Confusion Matrix?

Confusion matrix related metrics (where N is $TP + TN + FP + FN$)

Measure	Formula
error	$(fp + fn)/N$
accuracy = $1 - \text{error}$	$(tp + tn)/N$
tp-rate (hit rate)	tp/p
fp-rate (false alarm rate)	fp/n
precision	tp/p'
recall = tp-rate	tp/p
F-score	$\frac{\text{precision} \times \text{recall}}{(\text{precision} + \text{recall})/2}$
sensitivity = tp-rate	tp/p
specificity = $1 - \text{fp-rate}$	tn/n

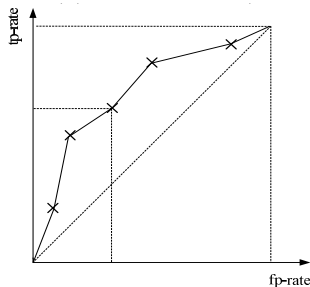
**Q. How do you measure how good your classification algorithm is?
You have unbalanced numbers, with our class being much larger or smaller than others**

You use the Matthews correlation coefficient. The number it yields is between plus or minus one. Plus one means perfect prediction, zero means no better than random.

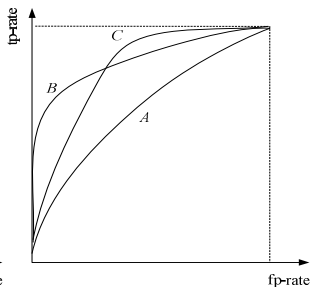
$$\frac{TP \times TN - FP \times FN}{\sqrt{(TP+FP)(TP+FN)(TN+FP)(TN+FN)}}$$

Q. ROC and AUC? Also given an example

E.g., suppose there is a threshold (or parameter) in your classification algorithm that determines whether you have an apple or a non apple object:



(a) Example ROC curve



(b) Different ROC curves for different classifiers

Process (how to split the data)

Q. Process description (Summary)?

Run fivefold cross-validation and select the best models.

Re-fit models on training set, evaluate them on the validation set, pick the best one.

Evaluate best model on the test set to gauge its ability to generalize to unseen data.

**What do we mean by training, best fit,
performance? Any risks?**

Q. But, what is training in ML anyway?

Most ML algorithms need to be trained. That is, you give them data and they look for patterns, or best fits, etc.

They know they are doing well when perhaps a loss function has been minimized, or the rewards have been maximized.

Q. Best fits? What do you mean by best fits? Performance?

First of all, when we say “fitting” in ML we are talking about model fitting.

Model fitting is a measure of how well a machine learning model generalizes¹⁸ to similar data on which it was *trained*.

A model that is well-fitted produces more accurate outcomes¹⁹.

¹⁸Yes, we are talking about ML model performance.

¹⁹An overfitted model matches the data too closely. An underfitted model doesn't match closely enough.

Q. How do you measure the performance of a ML model?

We do that by using a **cost function** or a loss function.

A cost function is used to represent how far away our model is from the real data.

A common way to do this is via the quadratic cost function²⁰:

$$J(\theta) = \frac{1}{2N} \sum_{n=1}^N (h_{\theta}(x^{(n)}) - y^{(n)})^2.$$

We are interested in the parameters that minimize this quadratic cost function.

This is called ordinary least squares (OLS).

One adjusts the mathematical model by varying parameters within the model, so as to *minimize the cost function*: giving the best model that fits the data.

²⁰This is just the sum of the squares of the vertical distances between the points and the straight line

Q. What are the risks of not splitting the data?

Overfitting or underfitting to the data

Inaccurate representation of how the model will generalize

Q. Any caveats on training and testing?

Over many epochs, if the test error begins to rise (and it's much bigger than the training error) then you have overfitted.

(Please refer to the Measuring success section to discuss ways in which we can void overfitting.)

Model optimization

Q. Model optimization outline?

We'll discuss the bias-variance trade-off.

We'll cover what we mean by bias and variance from a conceptual level.

Q. Bias and Variance in Machine learning?

Bias²¹, in ML, is the algorithm's tendency to consistently learn the wrong thing by not taking into account all the information in the data. (results in inaccurate predictions).

Variance²² is an algorithm's sensitivity to small fluctuations in the training data set.

²¹High bias is the result of the algorithm missing the relevant relations between features and target outputs

²²High variance is a result of the algorithm fitting to random noise in the training data

Q. Can you be more specific about Bias and Variance?

Bias is how far away (or error) the trained model is from the correct result *on average*. Where “on average” means over many goes at training the model, using different data:

$$\text{Bias}(\hat{f}(x')) = \underbrace{\mathbb{E}[\hat{f}(x)]}_{\text{Average error}} - f(x')$$

Variance is a measure of the magnitude of that error.

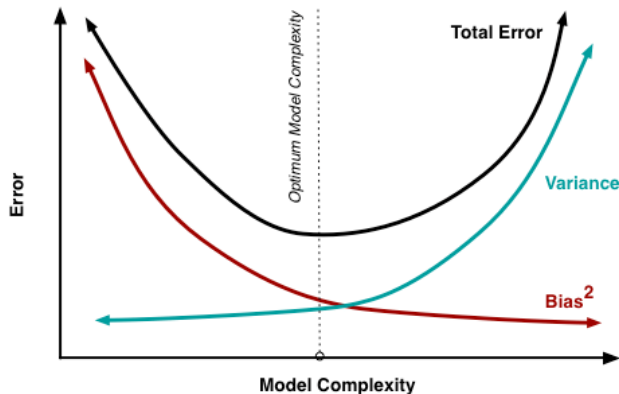
$$\text{Var}(\hat{f}(x')) = \mathbb{E}[\hat{f}(x)^2] - \mathbb{E}[\hat{f}(x')]^2$$

Bias and Variance tradeoff: As one is reduced, the other is increased²³.

²³This the matter of over-and-underfitting

Q. Bias and Variance tradeoff?

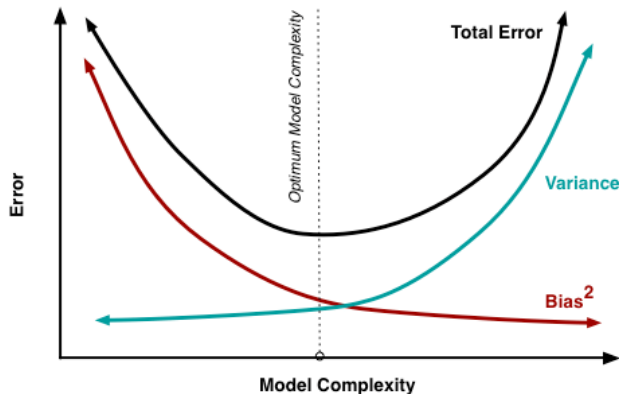
Total error = (Bias + Variance) + Irreducible Error



Model complexity is across the x axis and model error across the y axis. More complexity means higher variance. Lesser complexity means higher bias.

Q. Bias and Variance tradeoff?

Total error = (Bias + Variance) + Irreducible Error

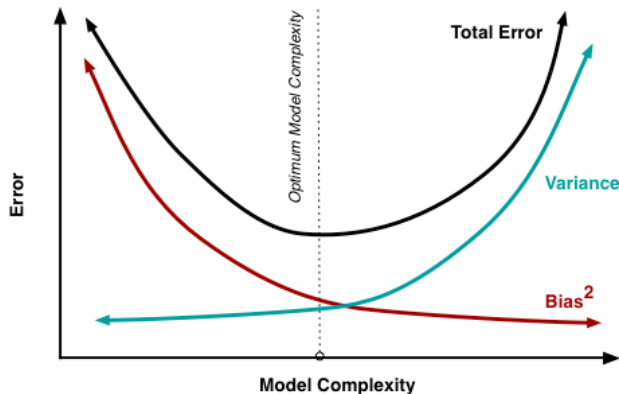


Model complexity is across the x axis and model error across the y axis. More complexity means higher variance. Lesser complexity means higher bias.

So this what bias/variance tradeoff is all about: **finding the right model complexity** that minimizes both bias and variance (I mean the total error as much as possible).

Q. Bias and Variance tradeoff?

Total error = (Bias + Variance) + Irreducible Error



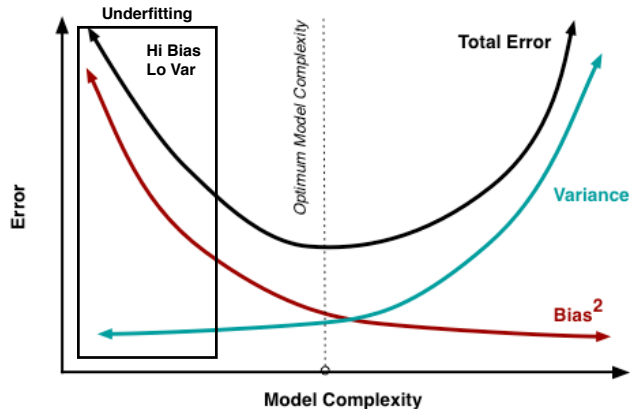
Model complexity is across the x axis and model error across the y axis. More complexity means higher variance. Lesser complexity means higher bias.

So this what bias/variance tradeoff is all about: **finding the right model complexity** that minimizes both bias and variance (I mean the total error as much as possible).

Total error is very high for very simple models and a very complex model, and then it bottoms out in the middle.

Q. What is Underfitting?

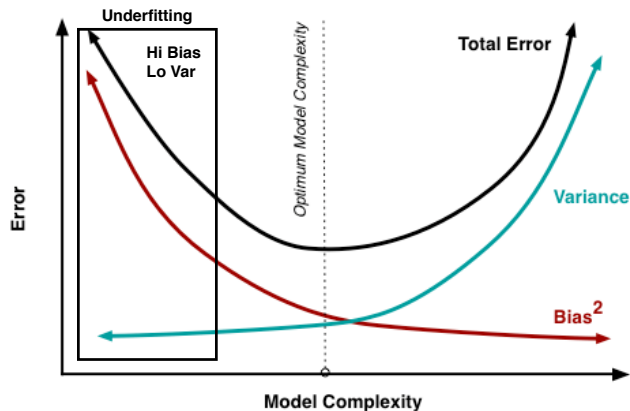
Underfitting occurs when an algorithm cannot capture the underlying trend of the data.



Underfitting happens when the model is too simple with high bias and low variance, and results in high total error.

Q. What is Underfitting?

Underfitting occurs when an algorithm cannot capture the underlying trend of the data.

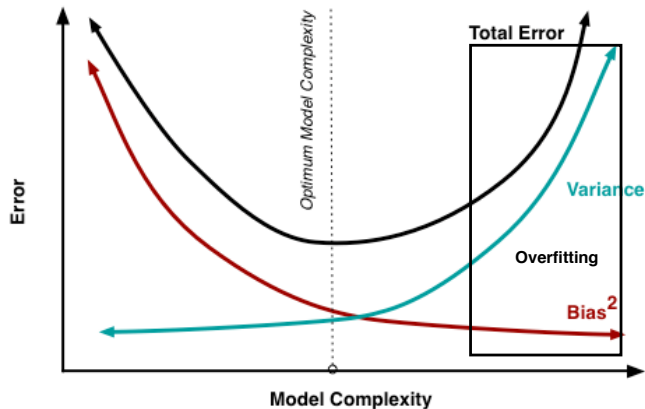


Underfitting happens when the model is too simple with high bias and low variance, and results in high total error.

Underfitting: High bias + low variance

Q. What is Overfitting?

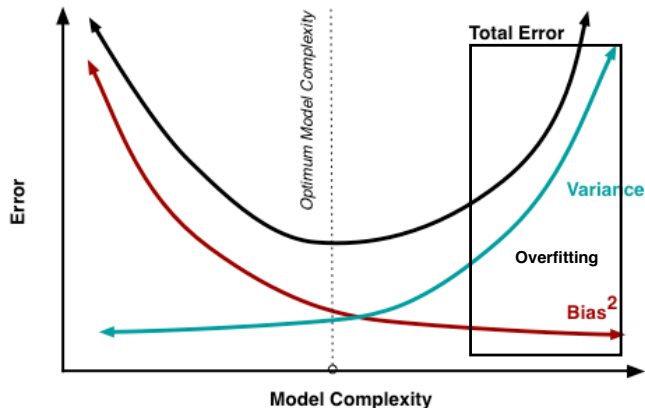
Overfitting occurs when an algorithm fits too closely to a limited set of data (training set?).



In other words, the model might just memorize the examples that it has seen in the training data.

Q. What is Overfitting?

Overfitting occurs when an algorithm fits too closely to a limited set of data (training set?).



In other words, the model might just memorize the examples that it has seen in the training data.

Overfitting: Low bias + high variance.

The actual process?

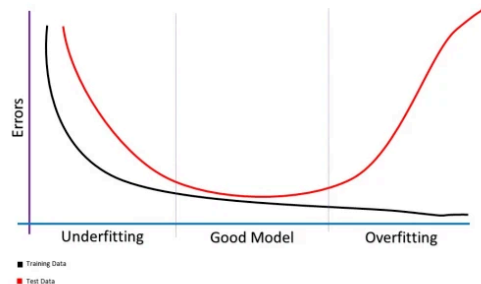
Q. How do you find the optimal tradeoff?

The goal is to find something in the middle²⁴ (a model with medium complexity); i.e.,

Optimal tradeoff: Low bias + Low variance.

OKAY, but how do you identify underfit and overfit?

With underfit, we'll have high training error and high test error.



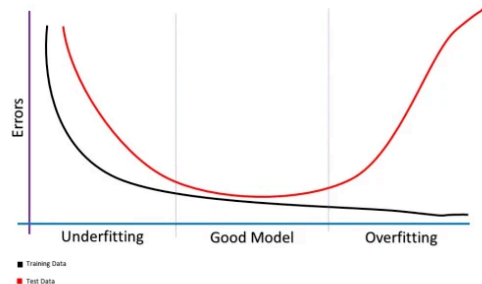
²⁴This would learn the true pattern in the data w/o memorizing every example in the training data.

Q. How do you find the optimal tradeoff?

The goal is to find something in the middle²⁴ (a model with medium complexity); i.e.,

Optimal tradeoff: Low bias + Low variance.

OKAY, but how do you identify underfit and overfit?



With underfit, we'll have high training error and high test error.

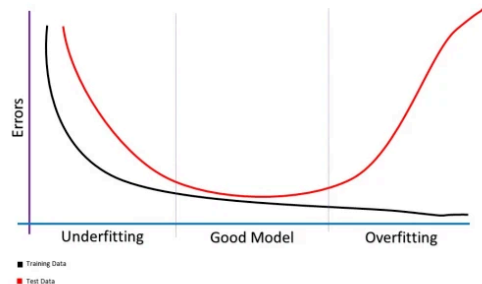
Optimal tradeoff, we'll have low training error and low test error.

²⁴This would learn the true pattern in the data w/o memorizing every example in the training data.

Q. How do you find the optimal tradeoff?

The goal is to find something in the middle²⁴ (a model with medium complexity); i.e.,
Optimal tradeoff: Low bias + Low variance.

OKAY, but how do you identify underfit and overfit?



With underfit, we'll have high training error and high test error.

Optimal tradeoff, we'll have low training error and low test error.

With overfit, we'll have low training error and high test error.

²⁴This would learn the true pattern in the data w/o memorizing every example in the training data.

Q. How do you tune a model for optimal complexity?

There are two methods to tune a model for optimal complexity:

Hyper-parameter tuning – choosing a set of optimal hyper-parameters for fitting a ML algorithm (e.g., linear regression)

Regularization – a technique used specifically to reduce overfitting by discouraging overly complex models in some way.

Q. What is a hyper-parameter?

A model **parameter** is a configuration variable that is internal to the model and *whose value can be estimated from data*.

A model **hyper-parameter** is a configuration that is external to the model and *whose value cannot be estimated from data, and whose value guides how the algorithm learns parameter values from the data*. E.g., depth of a decision tree is a model hyper-parameter vs ticket price or ticket class are model parameters.

Q. What is Regularization?

Regularization is a form of regression, which constrains/regularizes or shrinks the (learned) coefficient estimates towards zero. In other words, this technique reduces overfitting by discouraging learning a more complex model in some way.

The goal of Regularization is to allow enough flexibility for the algorithm to learn the underlying patterns in the data but provides guardrails so it doesn't overfit. See Occam's razor – whenever possible, choose the simplest model to a problem.

Q. Can you provide Regularization examples?

Ridge regression and lasso regression: adding a penalty to the loss function (See Model fitting slide) to constrain coefficients.

Dropout: some nodes are ignored during training which forces the other nodes to take on more or less responsibility for the input-output.

Epoch vs Batch size vs Iteration

Q. What is the difference between these terms?

To find out the difference between these terms you need to know some of the machine learning terms like Gradient Descent to help you better understand.

Q. What is Gradient Descent?

Gradient descent is an *iterative* optimization algorithm used to minimize some *convex* function by *iteratively* moving in the direction of steepest descent as defined by the negative of the gradient.

In ML, we use gradient descent to update the parameters of a ML model.

Start with an initial guess for each parameter θ_k . Then move θ_k in the direction of the slope.

Update all θ_k parameters simultaneously (known as batch gradient descent), and repeat until convergence.

Q. Unpacking Gradient Descent?

Gradient means the rate of inclination or declination of a slope.

Descent means we are dealing with the inclination of a slope.

The algorithm is iterative means that we need to get the results **multiple times** to get the most optimal result (minima of a curve).

Q. What is Stochastic Gradient Descent?

Similar to batch gradient descent except that you only update using one of the data points each time.

And that data point is chosen randomly. Hence the stochastic.

In other words, stochastic gradient descent pick an n at random and then update using one of the data points. Repeat, picking another data point at random, etc.

Q. What is an epoch?

In some ML methods, one uses the same training data many times, as the algorithm gradually converges, for example, in stochastic gradient descent. Each time the whole training set of data is used in the training that is called an **epoch**²⁵.

One might see a decreasing error as the number of epochs increases. But that doesn't mean your algorithm is getting better, it could easily mean that you are overfitting²⁶. To test for this you introduce a test data set, the data that you've held back.

²⁵Typically you won't get decent results until convergence after many epochs

²⁶This could happen if the learning algorithm has seen the training data so many times or epochs

Q. So, what is the right numbers of epochs?

Unfortunately, there is no right answer to this question. The answer is different for different data sets but you can say that the numbers of epochs is related to how diverse your data is. For example, do you have only black cats in your data set or is it much more diverse dataset?

Q. What is a batch?

As I said, you can't pass the entire data set into a ML algorithm at once. So, you divide data set into a number of batches or sets or parts.

Q. What is an iteration?

An iteration is the number of batches needed to complete one epoch.

Let's say we have 2000 training examples that we are going to use.

We can divide the dataset of 2000 examples into batches of 500 then it will take 4 iterations to complete 1 epoch.

End-to-End Pipeline

(starting from fit initial model)

Run fivefold cross-validation and select the best models.

Q. Fit a basic model using cross validation?

Goal: To understand what the baseline performance looks like ²⁷.

From sklearn, import RandomForestClassifier, and cross-val-score to compute the accuracy of the baseline model.

Next, we will select other models in order to beat our baseline.

²⁷We'll use only the training set

Q. Tuning Hyper-parameters?

Run grid search to find the optimal hyper-parameter settings for our models.

Our goal of this step is to find the optimal model that beats that baseline performance.

We use GridSearchCV²⁸, which yields multiple combinations of hyper-params values, to find the combinations that improves performance²⁹.

²⁸a wrapper around cross fail score that allows us to run grid search within cross validation.

²⁹Recall that a model is a configuration of the ML algorithm on specific params values

Re-fit models on full training set,
evaluate those models on the validation
set and pick the best one.

Q. Evaluating results on Validation set?

Now that we've done some hyper parameter tuning, and we have a good idea of what the best hyper parameter combinations are, let's evaluate these models on the validation set³⁰.

The process goes like this:

- 1 Look at additional performance metrics beyond just accuracy (e.g., precision and recall) and also at different ML algorithms.
- 2 Take the 3 best performing models and **re-fit** them using the whole training data³¹.
- 3 Evaluate these models on the validation set³².
- 4 Select the best performing model.

³⁰Now the performance shouldn't deviate too much from the performance we saw with the cross-validation.

³¹Why do we need to do that? Because these models were fit on only 80% of the training data; we need now the entire training data.

³²This is the truth test. If they are overfit or underfit, then they will fail here.

Evaluate that best model on the test set
to gauge its ability to generalize to
unseen data.

Q. Final model selection and evaluation on test set?

Next, we'll evaluate the best model on the test set. This will give us a truly unbiased view of how it should perform moving forward.

Why do we need to evaluate it on test data, given they've already evaluated on unseen data?

For the validation set was created by **randomly** distributing examples from the full data set. So if that validation set was slightly different perhaps we would have chosen a different model.

So by nature of this fact that testing on the validation set impacted what model was chosen as our final model, it's technically part of the model training process.

So this final test set is purely for evaluation purposes to see that it matches the performance that we've seen before and to give us more confidence in the performance of the model moving forward.

Machine Learning Algorithms

Linear Regression algorithm (Supervised learning)

Q. What is Regression?

Regression is a statistical process for estimating the relationship among variables, often used to make a prediction about some outcome.

From a Machine Learning perspective, the modeling of this relationship is done to ensure generalization – giving the model the ability to predict outputs (dependent variable) for inputs (independent variables) it has never seen before.

Q. What is Linear Regression?

Linear regression is one type of regression that is used when trying to predict a continuous target variable (output).

The way it works it attempts to model the relationship between two variables by fitting a linear equation (i.e., a line $Y = a + bX$) to observed data D , where X is the explanatory variable and Y is the dependent variable. For example, a modeler might want to relate the weights of individuals to their heights using a linear regression model.

To fit this regression line, it often uses **least-squares method**³³, which calculates the best-fitting line for D by minimizing the sum of the squares of the vertical deviations from each data point to the line (if a point lies on the fitted line exactly, then its vertical deviation is 0).

This minimization of error is wrt to values of a and b , which are found using **gradient descent** (training is basically finding these values.)

³³Because the deviations are first squared, then summed, there are no cancellations between positive and negative values.

Q. Linear regression on many dimensions?

We now have M independent, explanatory, variables, representing the features, and we write all x_s as vectors. For each of N data points we'll have the independent variable $x^{(n)}$ (e.g., squared footage of a house, number of rage bays) and the independent variable $y^{(n)}$ (property value).

We fit the linear function $h_{\theta}(x) = \theta^T x$ to the y_s , where θ is the vector of the as-yet-unknown parameters.

The cost function remains the same as in the one dimension LR: the quadratic function. And the coefficients are learned using batch or stochastic gradient descent.

Logistic Regression algorithm (Supervised learning)

Q. What is logistic regression?

Logistic regression or logistic classifier is a form of regression where the target variable is binary (zero or one, or true or false). It takes the form: $\frac{1}{1+e^{-(mx+b)}}$

Example: Let's say you wanted to examine the relationship between your basketball shooting accuracy and the distance that you shoot from. More specifically, you want a model that takes in “distance from the basket” in feet and spits out the probability that you'll make the shot (1 for a make, 0 for a miss).

In this example, the equation $-(mx + b)$, labeled as z , represents *log(odds of making shot)*³⁴, so the probability of making a shot (i.e., y) is $\frac{1}{1+e^{-z}}$. Here, m and b are the coefficients we're interested in finding thru optimization, and x is *distance from basket*.

³⁴odds = $P(\text{Event}) / [1-P(\text{Event})]$

Q. Logistic regression example?

Shooting Baskets. Generally, the further you get from the basket, the less accurately you shoot: when given a small distance, the model should predict a high probability and when given a large distance it should predict a low probability.

In logistic regression the output Y is in log odds. Let's say you shot 100 free throws and made 70. Based on this sample, your probability of making a free throw is 70%. Your odds of making a free throw is 2.33, which we want to bound between 0 and 1 using *log* and the *sigmoid* function as odds go from 0 to INF ³⁵:

We can write our logistic regression equation: $z = \log(mx + b)$, where m and b are the coefficients to be learned.

And to get probability from z , which is in log odds, we apply the sigmoid function:
 $\frac{1}{1+e^{-z}}$ ³⁶ In this case, a high probability means you'll be able to shoot and a low probability means you won't.

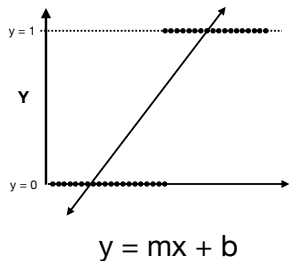
³⁵log odds are just probability stated another way

³⁶This shows how we can go from a linear estimate of log odds to a probability.

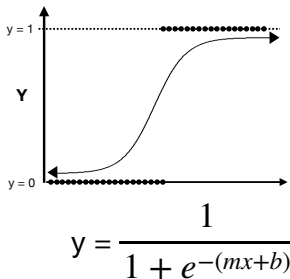
Q. Why won't Linear Regression work for binary target variables?

In other words, why do we need two regression algorithms?

Linear regression



Logistic regression

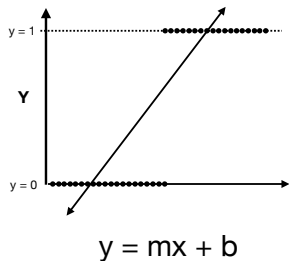


Linear regression for binary target variable will have a hard time try to come up with a best fit line that makes any sense.

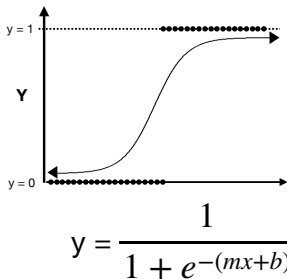
Q. Why won't Linear Regression work for binary target variables?

In other words, why do we need two regression algorithms?

Linear regression



Logistic regression

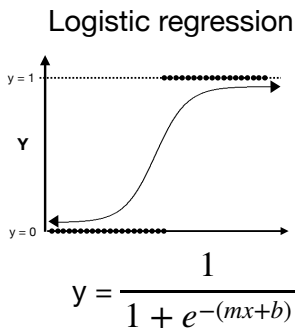
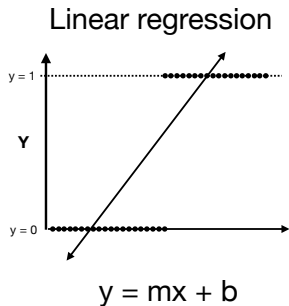


Linear regression for binary target variable will have a hard time try to come up with a best fit line that makes any sense.

It'll try to fit a line that fits all of the data and it will end up predicting negative values and values over one, which is impossible.

Q. Why won't Linear Regression work for binary target variables?

In other words, why do we need two regression algorithms?



Linear regression for binary target variable will have a hard time try to come up with a best fit line that makes any sense.

It'll try to fit a line that fits all of the data and it will end up predicting negative values and values over one, which is impossible.

Logistic regression is built off of a logistic or sigmoid curve (S shape) and will always be between zero and one, which makes it a better fit for a binary classification problem.

Q. When should you consider using Logistic Regression?

When to use? Consider using it any time you

- You have a binary target variable.
- You're interested in feature importance or having a better understanding of what's going on within the algorithm.
- You have well-behaved³⁷ data, need a **quick** initial benchmark.

When not to use? Consider not using it any time you

- You have a continuous target variable.
- You have a massive amount of data³⁸.
- You have unwieldy³⁹ data, performance⁴⁰ is the only thing that matters.

³⁷e.g., no many outliers, no many missing values, no complex relationships

³⁸e.g., lots of features and very few rows, or lots of rows and very few features

³⁹e.g., many outliers, many missing values, complex relationships

⁴⁰will usually do pretty well on any given problem, but will rarely be the best.

KMeans Clustering Algorithm

Q. What is KMeans?

KMeans is a unsupervised ML algorithm. It takes some data points as input and groups⁴¹ them into k clusters⁴² (the output). This k is called a hyper-parameter; a variable whose value we set before training.

The idea is simple: You have a bunch of vectors $\in \mathbb{R}^d$. **(1)** Then you init a bunch seed-guesses for the k centers (centroids) of the clusters. **(2)** You assign each vector to the nearest cluster. **(3)** And then, when everything is done, you calculate the mean values of the clusters (updated k centers). **(4)** And then you just do it again: you re-assign the nearest mean.

Then you keep on going until it converges⁴³.

⁴¹grouping is the training phase of KMeans, and uses the square of the L2 norm as its cost function

⁴²It uses the distance between points as a measure of similarity, based on k averages (i.e. means).

⁴³Once those centroids stop moving(no further change in cost), the algorithm stop

Q. When to use KMeans?

When to use it

You have unlabeled data and don't know the number of clusters within it.

You have a decently large data set (less than 10K) with a smaller number of dimensions, these data are numeric, or continuous.

When not to use it

High dimensional data, or data of varying sizes and density.

Messy data with lots of outliers (as it can centroids can be dragged by outliers).

Naive Bayes Classifier

Q. What is Naive Bayes?

Naive Bayes is a classification technique (Generative model) based on **Bayes Theorem** with an assumption of independence⁴⁴ among predictors (features).

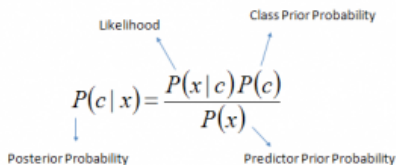
Some **benefits** include:

It's *easy to build* and particularly *useful for very large data sets*.

Along with *simplicity*, it's known to outperform even highly sophisticated classification methods (NB has better *resilience to missing* data than SVM).

⁴⁴In other words, it assumes that each feature makes an independent and equal contribution to the outcome.

Q. Bayes Theorem (Quick Overview)?



The diagram shows the Bayes' Theorem formula $P(c|x) = \frac{P(x|c)P(c)}{P(x)}$. Blue arrows point from labels to the corresponding parts of the formula: 'Likelihood' points to $P(x|c)$, 'Class Prior Probability' points to $P(c)$, 'Posterior Probability' points to $P(c|x)$, and 'Predictor Prior Probability' points to $P(x)$.

$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$$

Bayes theorem

Offers a way of calculating posterior probability $\Pr(c | x)$ from $\Pr(c)$, $\Pr(x)$, and $\Pr(x | c)$.

$\Pr(c | x)$ simply means, “given some feature vector $x_i \in X$, what is the probability of sample i belonging to class $c_j \in C$?”

The **objective function** of **Naive Bayes**: Maximize $\Pr(C | X)$ given the training data to formulate a decision rule for new data.

Q. How does Naive Bayes work?

First, this classifier uses the Bayes's theorem to calculate the probability of a class, say y_i , given a set of feature values (i.e. $\Pr(y_i \mid x_1, x_1, \dots, x_n)$) where the numerator ($\Pr(x_1, x_1, \dots, x_n \mid y_i)$) is probability of a specific combination of features given a class label and the denominator⁴⁵ ($\Pr(x_1) \Pr(x_2) \dots, \Pr(x_n)$) is the probability of all the evidence irrespective of our knowledge about the class label y_i . **And then** find the class label y_i with max probability.

More specifically:

1. Convert the data set into a frequency table (co-occurrence matrix)
2. Create Likelihood table by finding the probabilities like "Spam" probability (0.29) and probability of "No Spam" (0.64).
3. Now, use Naive Bayesian equation to calculate the posterior probability for each class. The class with the highest posterior prob. is the outcome of prediction.

⁴⁵A normalization term we often exclude.

Q. Example: Naive Bayes?

Problem: Players will play if weather is sunny. Is this statement is correct?

We can solve it using above discussed method of posterior probability.

$$\Pr(Y \mid \text{Sunny}) = \Pr(\text{Sunny} \mid Y) * \Pr(Y) / \Pr(\text{Sunny})$$

Here we have $\Pr(\text{Sunny} \mid Y) = 3/9 = 0.33$, $\Pr(Y) = 9/14 = 0.64$,
 $\Pr(\text{Sunny}) = 5/14 = 0.36$

Now, $\Pr(Y \mid \text{Sunny}) = 0.33 * 0.64 / 0.36 = 0.60$, which has higher probability.

Naive Bayes uses a similar method to predict the probability of different class based on various attributes. This algorithm is mostly used in text classification and with problems having multiple classes.

Perceptron

Perceptron is a single layer neural network (a multi-layer perceptron is called Neural Networks).

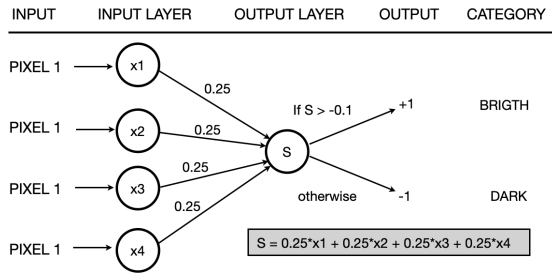
Q. What the Hell is Perceptron?

Perceptron is a linear binary classifier, used in supervised learning to classify data.

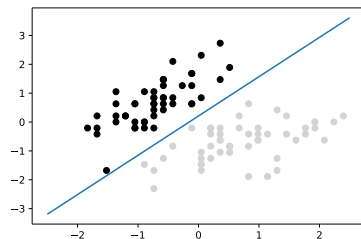
It consists of 4 components: (1) Input values or One input layer, (2) Weights and Bias, (3) Weighted sum, and (4) Activation function.

The perceptron works on these simple steps:

1. All inputs x are multiplied with their weights w .
2. Add all the multiplied values and call them Weighted Sum.
3. Apply that weighted sum to the correct Activation Function (e.g., rectified linear unit or simply ReLU) to return an output.



Q. When to use the Perceptron algorithm?



When to use it: A single perceptron can only be used to implement **linearly separable** functions. It takes both real and boolean inputs and associates some weights to them, along with a bias (the threshold of some sort). We learn the weights, then we get the function.

When not to use it: Non-linearly separable data.

Decision Trees Classifier

Q. Decision Trees?

Decision trees are a non-parametric model for supervised learning whose local region is identified in a sequence of recursive splits in a smaller number of steps.

It consists of internal decision nodes and terminal leaves. Each decision node implements a test function with discrete outcomes labeling the branches.

Q. How do Decision Trees work?

Decision trees (induction): Given an input, at each node, a test is applied and one of the branches is taken depending on the outcome.

This process starts at the root of the tree, and is repeated recursively until a leaf node is hit, at which point the value written in the leaf constitutes the output.

Q. Building a decision tree (a tree learning approach)?

Tree learning algorithms are greedy, and at each step, starting at the root with the complete training data, we look for the best split.

This splits the training data into two or n , depending on whether the chosen attribute is numeric or discrete (purer subsets⁴⁶).

We continue splitting (recursively) with the corresponding subset until no need to split anymore, at which point.

⁴⁶ → if all records in at that level belong to the same class.