

# Machine Learning: The Art and Science of Algorithms that Make Sense of Data

Tim Lawson

October 20, 2023

These notes are based on Flach 2012.

## Contents

5.1	Decision trees . . . . .	4
5.1.1	Purity of a leaf . . . . .	4
5.1.2	Impurity of a leaf . . . . .	4
5.1.3	Impurity of a tree . . . . .	5
5.1.4	Multi-class classification . . . . .	5
5.1.5	Purity and information gain . . . . .	5
5.2	Ranking and probability estimation trees . . . . .	6
5.3	Tree learning as variance reduction . . . . .	7
<b>8</b>	<b>Linear models</b>	<b>8</b>
8.1	The least-squares method . . . . .	8
8.2	The perceptron . . . . .	9
8.3	Support vector machines . . . . .	9
8.4	Obtaining probabilities from linear classifiers . . . . .	9
8.5	Notes . . . . .	9
<b>9</b>	<b>Distance-based models</b>	<b>11</b>
9.1	Distance metrics . . . . .	11
9.1.1	Examples . . . . .	11
9.2	Neighbours and exemplars . . . . .	12

## Theorems

5.1	Algorithm . . . . .	4
9.1	Definition (Metric) . . . . .	11
9.2	Definition (Pseudo-metric) . . . . .	11
9.3	Definition (Metric space) . . . . .	11
9.4	Definition ( $p$ -norm, $L_p$ norm) . . . . .	11
9.5	Definition (Minkowski distance) . . . . .	11
9.6	Definition (Manhattan distance) . . . . .	11
9.7	Definition (Euclidean distance) . . . . .	11
9.8	Definition (Chebyshev distance) . . . . .	12
9.9	Definition (0- "norm", $L_0$ "norm") . . . . .	12
9.10	Definition (Hamming distance) . . . . .	12
9.11	Definition (Mahalanobis distance) . . . . .	12
9.1	Theorem (Arithmetic mean minimises squared Euclidean distance)	12

**Trees** A tree is an undirected connected acyclic graph. A rooted tree is a tree in which a node is designated the root. The edges of a rooted tree may be directed away from or towards the root. The nodes of an  $m$ -ary tree have at most  $m$  children.

### Tree models

- A *tree model* is represented by a directed rooted tree. The branch nodes represent features and the leaf nodes represent instance space segments.
- A *feature tree* is represented by a binary directed rooted tree. There are two edges directed away from each branch node, each of which represents a mutually exclusive proposition about the value of the feature.
- A *decision tree* is represented by an  $m$ -ary directed rooted tree where  $m \geq 2$ . There are  $m_i$  edges directed away from each branch node  $i$ , each of which represents a possible value of the feature.

TODO:

- Disjunctive normal form
- Distributive equivalence  $A \vee (B \wedge C) \equiv (A \vee B) \wedge (A \vee C)$
- De Morgan laws  $\neg(A \vee B) \equiv \neg A \wedge \neg B$

**Expressivity** A decision tree represents a set of mutually exclusive logical expressions, which can be written in different equivalent forms. The expressions represented by a decision tree may not be equivalent to *conjunctive* expressions of individual features<sup>1</sup>. Because any logical expression may be written in disjunctive normal form, decision trees are maximally *expressive*, i.e., they can separate any data that is consistently labelled.

However, expressive hypothesis languages are prone to overfitting and one way to prevent overfitting is to choose a restrictive hypothesis language. Learning algorithms in expressive hypothesis spaces typically have an *inductive bias* towards simpler hypotheses, either implicitly by the search procedure or explicitly by a term in the loss function.

**Bias and variance** Low-bias models are more likely to overfit to the training data. Low-variance models change by a small amount when the training data changes by a small amount.

Low-variance, high-bias models are preferable when there is limited training data and overfitting is a concern. High-variance, low-bias models are preferable when there is plenty of training data but underfitting is a concern.

---

<sup>1</sup>But they may be equivalent to conjunctive expressions of *conjunctive features*. This is called *constructive induction*.

**Learning algorithms** A feature tree represents conjunctive concepts in the hypothesis space. The learning problem is to choose the best conjunctive concepts to solve a task.

Algorithm 5.1 (pseudo-code) is a generic learning procedure. It is a *divide-and-conquer* algorithm: it splits the data into subsets, learns a tree for each subset, and combines them. It is also a *greedy* algorithm: it always chooses the best feature values to split the data at a given step, which may be sub-optimal. An optimal but more computationally expensive alternative is to search for the best feature values to split the data over all steps.

**Algorithm 5.1.**

```
# Returns true if data can be given a single label.
def is_homogeneous(data)

# Returns the best label for data.
def label(data)

# Returns the best feature values to split data.
def find_feature_values(data, features)

# Returns the subsets of data for each feature value.
def find_split(data, feature, values)

def grow(tree, data, features):
    if is_homogeneous(data):
        tree.add_leaf(label(data))

    feature, values = find_feature_values(data, features)

    for subset in find_split(data, feature, values):
        if len(subset) > 0:
            grow(tree[feature], subset, features)
        else:
            tree[feature].add_leaf(label(subset))
```

## 5.1 Decision trees

### 5.1.1 Purity of a leaf

For a classification task, a set of instances is *homogeneous* if the instances belong to the same class. Therefore, **def** label returns the majority class. The *purity* of a set of instances is the proportion of instances that belong to the majority class. It is proportional to the *empirical probability*  $\hat{p}$ .

### 5.1.2 Impurity of a leaf

**def** find\_feature\_values returns the feature values that maximise the purity (minimise the impurity) of the subsets. In terms of  $\hat{p}$ , the impurity  $f$  must obey the following constraints:

- $f(\hat{p}) = 0 : \hat{p} \in 0, 1$ , i.e., it is zero if the subset is homogeneous

- $f(\dot{p}) = f(1 - \dot{p})$ , i.e., it is symmetric about  $\dot{p} = \frac{1}{2}$
- $\arg \max_{\dot{p}} f = \frac{1}{2}$ , i.e., it is maximal when  $\dot{p} = \frac{1}{2}$

Some examples of impurity functions are:

- *Minority class*  $\min(\dot{p}, 1 - \dot{p})$   
The error rate, i.e., the proportion of instances that are labelled incorrectly if they are labelled with the majority class.
- *Gini index*  $2\dot{p}(1 - \dot{p})$   
The expected error if we label instances randomly.
- *Entropy*  $-\dot{p} \log_2 \dot{p} - (1 - \dot{p}) \log_2 (1 - \dot{p})$   
The expected number of bits encoded by the class of a random instance.

### 5.1.3 Impurity of a tree

The impurity of a set of mutually exclusive leaves, i.e., a decision tree, is the weighted average of the impurities of the leaves:

$$f(\{D_i \mid i \in 1, \dots, n\}) = \frac{1}{|D|} \sum_{i=1}^n |D_i| f(\dot{p}_i) \quad (1)$$

For binary classification, we can find  $f(\{D_+, D_-\})$  from  $f(\dot{p}_+)$  and  $f(\dot{p}_-)$  geometrically: First, we draw a straight line between  $(\dot{p}_+, f(\dot{p}_+))$  and  $(\dot{p}_-, f(\dot{p}_-))$ . The line represents the possible weighted averages of  $f(\dot{p}_+)$  and  $f(\dot{p}_-)$ . Given that  $\dot{p} = \frac{|D_+|}{|D|} \dot{p}_+ + \frac{|D_-|}{|D|} \dot{p}_-$ ,  $f(\dot{p})$  is the point on the line that corresponds to  $\dot{p}$ .

### 5.1.4 Multi-class classification

Impurity functions can be generalized to  $k > 2$  classes, e.g.:

- *k-class Gini index*  $\sum_{i=1}^k \dot{p}_i (1 - \dot{p}_i)$
- *k-class entropy*  $\sum_{i=1}^k -\dot{p}_i \log_2 \dot{p}_i$

### 5.1.5 Purity and information gain

To split a parent node  $D$  into children  $\{D_i \mid i = 1, \dots, n\}$ , we typically choose the feature that maximises the *purity gain*:

$$f(D) - f(\{D_i \mid i = 1, \dots, n\}) \quad (2)$$

If  $f(D)$  is the entropy, this is called the *information gain*. It measures the increase in information about the class gained by including the feature. A ‘best split’ algorithm finds the feature that minimises  $f(\{D_i \mid i = 1, \dots, n\})$ .

## 5.2 Ranking and probability estimation trees

A grouping classifier can be used to rank instances by learning an ordering on its instance-space segments. Decision trees can access the class distributions (empirical probabilities) of the segments, from which an ordering can be derived that is optimal for the training data. This is not possible for some other grouping classifiers.

The ordering is optimal because it produces a convex ROC curve. The ROC curve is convex because its segments are sorted in decreasing order of slope. The slope of a segment is  $\frac{\dot{p}}{1-\dot{p}}$  and, because the slope is a monotonic function of  $\dot{p}$ , sorting the segments in decreasing order of  $\dot{p}$  is equivalent to sorting them in decreasing order of slope.

The empirical probability of a parent node is a weighted average of the empirical probabilities of its children (see 1):

$$\dot{p} = \frac{1}{|D|} \sum_{i=1}^n |D_i| \dot{p}_i \quad (3)$$

But this does not constrain the empirical probabilities of a parent's children, so we cannot find the ordering of segments from the tree structure.

TODO:

- Interpretation of splits in terms of coverage curves. To add a split: split the line segment of the ROC curve into  $k > 2$  segments and re-sort the segments in decreasing order of slope. Sorting the segments ensures that the ROC curve is convex.
- Turning a feature tree into a decision tree (classifier), ranking tree, or probability estimation tree.
  - Decision tree (classifier): choose the operating conditions and find the optimal point under those conditions.
  - Ranking tree: order the segments in decreasing order of empirical probability.
  - Probability estimation tree: predict the empirical probabilities of the segments (applying smoothing).
- Pruning trees.
  - Merging all leaves in a subtree
  - Only recommended for classification and when you can define the operating conditions
  - E.g., reduced-error pruning with a pruning set
  - Never improves accuracy over the training data
- Sensitivity to imbalanced classes.
  - Oversampling the minority class. Applies to any model without changing the model itself. But increases training time and may not change the model(!).
- Relative impurity.

- The relative impurity of a child node is its weighted impurity in proportion to its parent node's impurity.
- Some impurity measures are invariant with respect to the class distribution. E.g., the square root of the Gini index, which minimises the relative impurity.
- Impurity measures that vary with the class distribution produce splitting criteria that emphasise child nodes with more instances. E.g., the Gini index and entropy.

How to train a decision tree:

1. Prioritise ranking performance
2. Use an impurity measure that is invariant with respect to the class distribution. Otherwise, oversample the minority class to balance the class distribution.
3. Apply Laplace or add- $k$  smoothing to the empirical probabilities.
4. Given the operating conditions, select the best operating point on the ROC curve.
5. Optionally, prune subtrees whose leaves are homogeneous.

### 5.3 Tree learning as variance reduction

TODO:

- Adapting decision trees to regression and clustering tasks.
- Variance of a Bernoulli distribution.
- Overfitting, pruning, and model trees.
- Cluster and split dissimilarity.

## 8 Linear models

- Linear models are defined in terms of the geometry of the instance space.
- Real-valued features are not generally intrinsically geometric.
- However, we can use geometric concepts to structure the instance space (e.g., lines and planes) and represent similarity by distance.

Linear models are simple:

- They are parametric: they have a fixed structure that is defined by numeric parameters that are learned from the training data. By contrast, tree and rule models are non-parametric: their structure is not fixed prior to learning.
- They are stable (have low variance): small variations in the training data have a small effect on the learned model. Tree models have high variance.
- They are unlikely to overfit the training data because they have relatively few parameters (have high bias). However, they sometimes underfit the training data.

### 8.1 The least-squares method

The least-squares method can be used to learn linear models for classification and regression. It finds a function estimator that minimises the sum of squared residuals (differences between the actual and estimated values).

**Univariate linear regression** Let  $\{(x_i, y_i) \mid i \in 1..n\}$  be a set of instances. Approximate the true function  $f(x_i) = y_i$  by a linear function  $f'(x_i) = a + bx_i$ . Univariate linear regression finds  $a, b$  such that the sum of squared residuals  $\sum_{i=1}^n (y_i - (a + bx_i))^2$  is minimized.

When the sum of squared residuals is minimized, its partial derivatives with respect to  $a$  and  $b$  are zero:

$$\frac{\partial}{\partial a} \sum_{i=1}^n (y_i - (a + bx_i))^2 = -2 \sum_{i=1}^n (y_i - (a + bx_i)) = 0 \quad (4)$$

$$\frac{\partial}{\partial b} \sum_{i=1}^n (y_i - (a + bx_i))^2 = -2 \sum_{i=1}^n (y_i - (a + bx_i))x_i = 0 \quad (5)$$

TODO

- Translation does not affect the regression coefficient, only the intercept. We can zero-centre the  $x$ -values by subtracting the mean  $\bar{x}$ .
- If we normalize  $x$  to have unit variance, then the regression coefficient is the covariance between the normalized  $x$  and  $y$ .



The least-squares solution is equivalent to the maximum likelihood estimate given the assumptions that the true function is linear but normally-distributed noise is added to the instance  $y$ -values. If noise is added to only the  $y$ -values, then it is called *ordinary* least squares, which has a unique solution. If noise is added to both  $x$ - and  $y$ -values, then it is called *total* least squares, which does not necessarily have a unique solution.

Zero-centred matrix, scatter matrix, covariance matrix.

**Multivariate linear regression** Matrix form and homogeneous coordinates. Transformation to decorrelate, centre and normalize features. If the features are assumed to be uncorrelated, a multivariate linear regression problem decomposes into a set of univariate linear regression problems. It is computationally expensive to invert the scatter (covariance) matrix.

**Regularization** Least-squares regression can be unstable. Instability demonstrates a tendency to overfit. Regularization helps to avoid overfitting by constraining the weight vector.

- Shrinkage: makes the average magnitude of the weights small. This adds a scalar parameter to the diagonal of the scatter matrix, which improves the numerical stability of matrix inversion. Least-squares regression with shrinkage is called ridge regression.
- Lasso (least absolute shrinkage and selection operator): This adds the sum of the absolute weights ( $L_1$  regularization). This makes the magnitude of some weights smaller but sets others to zero, i.e., it favours sparse solutions.

## 8.2 The perceptron

## 8.3 Support vector machines

## 8.4 Obtaining probabilities from linear classifiers

## 8.5 Notes

When to favour models with different characteristics? E.g., the quantity and quality training data.

Normalization (zero-centre, unit variance) Write up the equivalencies between correlation coefficients etc.

Regularization Relation to, e.g., Bayesian priors Technically, e.g., sparsity (Occam's razor). Regularization changes the optimal solution (it's included in the loss)

Correlation (in the extreme case, two copies of the same feature) — your problem is underspecified, i.e., there are infinitely many solutions (which are combinations of the two). 'Spikiness of the fitness landscape' (high variance). Regularization: decreasing dependence on the data, i.e., increasing bias and decreasing variance.

When to choose ridge (L2) vs lasso (L1) regularization? An elastic net uses a weighted combination of the two where the weight is a hyperparameter that you can tune.

Why does lasso produce sparse solutions? With Euclidean distance, the set of points at equal distance is a circle. If you change the exponent, e.g., Minkowski at  $d = 3$ , that changes shape. E.g.  $d = 1$  ‘pulls you’ towards a solution on one of the axes, i.e., towards one or the other feature instead of a combination of both (i.e., sparsity). Vector field analysis.

LP-norm where  $p$  is an integer.

Multivariate.

## 9 Distance-based models

A distance-based model is generally comprised of:

- a distance metric (section 9.1);
- a set of exemplars (centroids or medoids); and
- a distance-based decision rule.

### 9.1 Distance metrics

**Definition 9.1** (Metric). *A metric is a function  $d : M \times M \rightarrow \mathbb{R}$ , where  $M$  is a set of points, such that:*

1.  $d(x, x) = 0 \ \forall x \in M$  (the distance from a point to itself is zero)
2.  $d(x, y) > 0 \ \forall x, y \in M, x \neq y$  (positivity)
3.  $d(x, y) = d(y, x) \ \forall x, y \in M$  (symmetry)
4.  $d(x, z) \leq d(x, y) + d(y, z) \ \forall x, y, z \in M$  (triangle inequality)

**Definition 9.2** (Pseudo-metric). *A pseudo-metric is a metric where the condition of positivity is replaced by non-negativity, i.e.,  $d(x, y) \geq 0 \ \forall x, y \in M$ .*

**Definition 9.3** (Metric space). *A metric space is an ordered pair  $(M, d)$  where  $M$  is a set of points and  $d$  is a metric on  $M$ .*

#### 9.1.1 Examples

**Definition 9.4** ( $p$ -norm,  $L_p$  norm). *The  $p$ -norm of a vector  $\vec{x} \in \mathbb{R}^n$  is:*

$$\|\vec{x}\|_p = \left( \sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}} \quad (6)$$

**Definition 9.5** (Minkowski distance). *The Minkowski distance of order  $p \in \mathbb{N}_1$  between two vectors  $\vec{x}, \vec{y} \in \mathbb{R}^n$  is:*

$$D_p(\vec{x}, \vec{y}) = \left( \sum_{i=1}^n |x_i - y_i|^p \right)^{\frac{1}{p}} = \|\vec{x} - \vec{y}\|_p \quad (7)$$

**Definition 9.6** (Manhattan distance). *The Manhattan distance between two vectors  $\vec{x}, \vec{y} \in \mathbb{R}^n$  is the Minkowski distance of order  $p = 1$ :*

$$D_1(\vec{x}, \vec{y}) = \sum_{i=1}^n |x_i - y_i| \quad (8)$$

**Definition 9.7** (Euclidean distance). *The Euclidean distance between two vectors  $\vec{x}, \vec{y} \in \mathbb{R}^n$  is the Minkowski distance of order  $p = 2$ :*

$$D_2(\vec{x}, \vec{y}) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \quad (9)$$

**Definition 9.8** (Chebyshev distance). *The Chebyshev distance between two vectors  $\vec{x}, \vec{y} \in \mathbb{R}^n$  is the Minkowski distance of order  $p \rightarrow \infty$ :*

$$D_\infty(\vec{x}, \vec{y}) = \lim_{p \rightarrow \infty} \left( \sum_{i=1}^n |x_i - y_i|^p \right)^{\frac{1}{p}} = \max_{i=1}^n |x_i - y_i| \quad (10)$$

Minkowski distances are translationally invariant but not scale-invariant. Euclidean distance is the only Minkowski distance that is rotationally invariant.

**Definition 9.9** (0-“norm”,  $L_0$  “norm”). *The 0-“norm” of a vector  $\vec{x} \in \mathbb{R}^n$  is the number of non-zero elements in  $\vec{x}$ :*

$$\|\vec{x}\|_0 = \sum_{i=1}^n |x_i|^0 \quad (11)$$

The 0-“norm” is not a norm because it is not homogeneous.

**Definition 9.10** (Hamming distance). *The Hamming distance between two binary strings  $\vec{x}, \vec{y}$  of length  $n$  is the number of bits in which they differ:*

$$D_0(\vec{x}, \vec{y}) = \sum_{i=1}^n |x_i - y_i|^0 = \sum_{i=1}^n \mathbb{I}(x_i \neq y_i) \quad (12)$$

The edit or *Levenshtein distance* generalises the Hamming distance to non-binary strings of different lengths.

**Definition 9.11** (Mahalanobis distance). *The Mahalanobis distance between two vectors  $\vec{x}, \vec{y} \in \mathbb{R}^n$ , where  $\Sigma$  is the covariance matrix, is:*

$$D_M(\vec{x}, \vec{y} \mid \Sigma) = \sqrt{(\vec{x} - \vec{y})^\top \Sigma^{-1} (\vec{x} - \vec{y})} \quad (13)$$

Euclidean distance is the Mahalanobis distance where the covariance matrix is the identity matrix.

## 9.2 Neighbours and exemplars

**Theorem 9.1** (Arithmetic mean minimises squared Euclidean distance). *The arithmetic mean  $\vec{\mu}$  of a set of points  $X \in \mathbb{R}^n$  is the point with the minimum sum of squared Euclidean distances to the points in  $X$ :*

$$\arg \min_{\vec{y}} \sum_{\vec{x} \in X} \|\vec{x} - \vec{y}\|_2^2 = \vec{\mu} \quad (14)$$

*Proof.* The gradient of the sum of squared Euclidean distances is:

$$\begin{aligned} \nabla_{\vec{y}} \sum_{\vec{x} \in X} \|\vec{x} - \vec{y}\|_2^2 &= -2 \sum_{\vec{x} \in X} (\vec{x} - \vec{y}) \\ &= -2 \sum_{\vec{x} \in X} \vec{x} + 2|X|\vec{y} \end{aligned}$$

If the gradient is the zero vector, then:

$$\vec{y} = \frac{1}{|X|} \sum_{\vec{x} \in X} \vec{x} = \vec{\mu}$$

□

Minimising the sum of squared Euclidean distances is equivalent to minimising the *average* squared Euclidean distance. The *geometric median* minimises the sum of Euclidean distances. However, there is no closed-form expression for the geometric median of multivariate data.

## References

Flach, Peter (Sept. 2012). *Machine Learning: The Art and Science of Algorithms That Make Sense of Data*. 1st ed. Cambridge University Press.