

## Introduction to Machine Learning. Part 1: Problems and Approaches

### Machine Learning Problems

**Machine Learning Problem Description.** The most general description of Machine Learning is something like this:

*Given a number of observations, predict what other observations would look like.*

A somewhat broader definition states simply *find insight in data*.

There are a number of shapes that Data Mining/Machine Learning can take. Each "shape" or problem type is characterized by some unique circumstances/needs/questions, but they are also interconnected. One of our key goals for the course is to understand these interconnections.

**Machine Learning/Data Mining problems** Aggrawal outlines four categories of Data Mining problems:

- **Pattern analysis:** search of interesting patterns in data
- **Classification:** prediction of a value of a (categorical) feature associated with each data point
- **Clustering:** partitioning of objects in a dataset into groups such that objects each group are more similar to each other than to objects in other groups.
- **Outlier analysis:** finding data points in a dataset that exhibit anomalous or unusual values

In my view, this breakdown needs to be expanded and adjusted. In CSC 466 and in DATA 401 I usually give the following nomenclature of Machine Learning problems:

- **Regression problems:** prediction of a value of a numeric feature associated with each data point
- **Classification problems:** prediction of a value of a categorical feature associated with each data point
- **Clustering problems:** partitioning of objects in a dataset into groups such that objects each group are more similar to each other than to objects in other groups.
- **Recommendation problems:** prediction of a value of a (numeric or categorical) feature in a sparse matrix.

Another important issue here is to distinguish between *problems* and *approaches* that carry the same name. **Problems** are questions and challenges. **Approaches** are solutions. For example:

- *Outlier detection problem* can be solved using classification algorithms
- *Outlier detection problem* can be solved using clustering algorithms
- *Classification problem* can be solved using pattern analysis algorithms
- *Recommendation problem* can be solved using regression algorithms

and so on.

Here is another attempt at approaching this.

**Type of insight: Local or Global:** What is the scope of the problem w.r.t. the dataset?

- **Local.** Some ML/DM problems require local insight, i.e., analysis of a *part of the dataset*. While such analysis may involve looking at the entirety of the data available, the output of the analysis is **local** - i.e., it concerns only a subset of the data. Machine Learning problems with local scope are:
  - **Outlier Detection:** only data points that do not "fit" the rest of the data are of interest to us.
  - **Pattern Analysis:** a pattern is *explicitly* a subset of data in which a specific dependency/behavior is observed.
  - **Recommendations:** in some scenarios, the problem of giving a recommendation involves a single data point and a subset of the entire dataset. Global-scope versions of this problem also exist.
- **Global.** These types of analysis target the entire dataset will provide insight about every single data point in it. Often, when people say "Machine learning" or "data mining" they mean these kinds of problems.
  - **Regression:** regression predicts the value of a target variable for each point in the dataset as well as for any unseen data point.

- **Classification:** classification predicts the class for each data point in the dataset as well as for any unseen data point.
- **Clustering:** the entire dataset is partitioned, the status of each data point (i.e., what cluster it belongs to, if any) is of interest.
- **Recommendations:** global scenarios involve predicting an expected score for a specific item for all data points (users) in a dataset, and building predictive models for unseen data.

**What is the final goal of the analysis?** There are multiple final goals, and the same approaches can be used to meet different goals.

- **Prediction.** Given an unseen data point predict its properties. This is a common task for **classification**, **regression** and **recommendation** problems, as well as **outlier detection**. In some contexts **clustering** can be viewed from a predictive point of view as well.
- **Description.** Describe the insight gained in the collected data. **Clustering** and **pattern analysis** concentrate on insight from the collected data. Some **Outlier detection** scenarios are descriptive.
- **Interpretation.** The bread-and-butter of statistical analysis: we are interested in explanation of the reasons *why* the data behaves the way it does. **Classification** and **regression** are often used for *interpretation* rather than *prediction*, and these two uses of these approaches differ. Explanatory models can be constructed for **recommendations**, **outlier detection** and **pattern analysis** as well.

**What data is available?** This addresses the classic *supervised* vs. *unsupervised* learning split.

- **Supervised Learning methods.** Supervised learning relies on **ground truth**, i.e., known values of the variable(s) of interest in the provided dataset (often referred to as training data) to infer predictions or interpretations. Supervised learning problems are
  - **Regression:** training data has values for the numeric target variable
  - **Classification:** training data has class assignment
  - **Recommendation:** training data has item scores for a subset of customers (some people may refer to such situations as "semi-supervised learning")
  - **Outlier detection:** when viewed as a classification problem, the training set contains data points labelled as outliers/anomalies.
- **Unsupervised Learning methods.** Unsupervised learning relies on gaining insight in the absence of ground truth.
  - **Clustering:** can be thought of as a classification problem when the class variable is unknown.

- **Pattern Analysis:** examples of patterns of interest are typically NOT provided. Pattern analysis methods *can be used* in supervised learning settings though.
- **Outlier detection:** when viewed as a clustering problem, the available dataset may not include any outlier designations by itself.

Before proceeding with more detailed discussion of different types of machine learning problems, we first must examine the notion of *observations*.

## Data for Machine Learning: Objects, Features, Targets

In most cases, data for machine learning problems is a collection of *object descriptions*, where each *object* is described as a *collection of features*.

**Features.** A *feature* is an individual measurable property of an object or phenomenon that is being observed. A *feature* has the following properties:

- **Name.** An unique identifier that distinguishes the feature from all other features considered as part of the same machine learning problem.
- **Domain.** A set of values the feature can take (i.e., the set of possible measurements that can be made of the feature). Feature ranges can be finite or infinite.

**Types of features.** Depending on the specific range of feature values, features can be classified into the following categories:

1. **Numeric.** Numeric features have domains that are infinite or finite (but usually rather large) sets of numbers.
2. **Categorical.** Categorical features are features with domains that are not numeric. Among categorical features, we identify
  - (a) **Nominal features.** These are features whose possible values admit no order.
  - (b) **Ordinal features.** These features have possible values that can be meaningfully ordered, and thus, can be represented as numbers  $\{1, \dots, N\}$  where  $N$  is the size of the domain of the feature.

**Example.** Temperature of a hospital patient measured in degrees Celcius is a *numeric feature* whose domain is the range  $[30.0, 45.0]$ .

The patient's **native language** is a *categorical nominal feature*. Its domain may be rather large (total number of languages in the world is around 6900), by the possible values of this feature, e.g., "English", "Spanish", "Mandarin", "Russian", etc. do not have an ascribed order.

The patient's **level of pain** is a *categorical ordinal feature*. Its domain may be  $\{\text{no, slight, moderate, severe}\}$ . There is a natural order that can be imposed on

these values:  $\text{no} < \text{slight} < \text{moderate} < \text{severe}$ , which makes this feature *ordinal*. The domain can be represented as  $\{1, 2, 3, 4\}$  (or  $\{0, 1, 2, 3\}$ ) with numbers retaining their property of order.

**Synonyms.** We use the terms *feature*, *attribute*, and *variable* (or *independent variable*) interchangeably.

**Observations.** In *statistics*, an **observation** is a measured value, at a particular moment of time, of a specific feature.

In this class, we use the term **observation** in a more general way to refer to a collection of measured values of a particular set of features, that *describe a specific object or phenomenon*.

**Formalizing.** A given machine learning problem describes a set of objects or a set of phenomena by establishing a *set of features*  $\mathcal{A} = \{A_1, \dots, A_n\}$  whose values combined produce a *complete* (from the perspective of the specific machine learning problem) description of a single object/single phenomenon. The values for each feature  $A_i, i = 1 \dots n$  come from the set  $D_i = \text{dom}(A_i)$ .

A vector  $\mathbf{x} = (x_1, \dots, x_n)$  of values, where  $(\forall i = 1 \dots n)(x_i \in \text{dom}(A_i))$  is called a **data point**. Often, without loss of generality, we refer to data points as *objects*, *phenomena*, *entities*, *points*, *records*, *tuples*, *feature-vectors* or use a domain-specific name to identify them.

**Example.** Consider a collection of features  $\{\text{Name}, \text{Language}, \text{Temperature}, \text{PainLevel}\}$  representing the name of a hospital patient, their native language, their temperature at admission time, and their self-reported pain level.

A description of a single patient may be a vector of values

$$\mathbf{x} = (\text{"Mary Smith"}, \text{"English"}, 38.2, \text{"moderate"}).$$

We can refer to  $\mathbf{x}$  in a number of ways:

- *vector of feature values*, *feature-vector*, or simply *vector*
- *data point*
- *observation*
- *object* (although in this case it is a somewhat awkward term)
- *record* or *tuple*
- *patient record* (this is a domain-specific term that carries knowledge of the semantics of the data)
- *patient*

**Dataset.** Let  $\mathcal{A} = \{A_1, \dots, A_n\}$  be a set of features. Let  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ , where

$$(\forall i = 1 \dots m)\mathbf{x}_i = (x_{i1}, \dots, x_{in}),$$

and

**Dataset.** Let  $\mathcal{A} = \{A_1, \dots, A_n\}$  be a set of features. Let  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ , where

$$(\forall i = 1 \dots m) \mathbf{x}_i = (x_{i1}, \dots, x_{in}),$$

and

$$(\forall j = 1 \dots n) x_{ij} \in \text{range}(A_j).$$

$X$  is called a *collection of datapoints*, or a *dataset*.

$X$  is called a *collection of datapoints*, or a *dataset*.

## Algebraic View of Datasets

**Dataset as a matrix.** A dataset  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$  can be viewed as a *matrix*  $X$  constructed as follows:

- The *rows* of  $X$  are vectors  $\mathbf{x}_1, \dots, \mathbf{x}_m$ .
- The *columns* of  $X$  are individual features  $A_1, \dots, A_n$ , with each column containing the values of a single feature from all data points of the dataset.

Dataset  $X$  can be visualized as a matrix as follows:

$$X = \left( \begin{array}{c|cccc} & A_1 & A_2 & \dots & A_n \\ \hline \mathbf{x}_1 & x_{11} & x_{12} & \dots & x_{1n} \\ \mathbf{x}_2 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_m & x_{m1} & x_{m2} & \dots & x_{mn} \end{array} \right)$$

**Size of a dataset.** The **size** of a dataset  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , denoted  $|X|$  is the number of data points in it, i.e.,

$$|X| = m$$

The **dimensionality** of the dataset  $X$ , denoted  $\text{dim}(X)$  is the length of the vectors in it. If

$$\mathbf{x}_i = (x_{i1}, \dots, x_{in}),$$

then

$$\text{dim}(X) = n.$$

**Numeric vector-spaces.** Let  $\mathcal{A} = \{A_1, \dots, A_n\}$  and let  $(\forall i \in 1 \dots n) \text{dom}(A_i) \subseteq \mathbb{R}$ . In this case, given a dataset  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  of points in the feature-space  $\mathcal{A}$ , we can write for each data point  $\mathbf{x}_i$  that

$$\mathbf{x}_i = (x_{i1}, \dots, x_{in}) \in \mathbb{R}^n,$$

i.e.,  $\mathbf{x}_i$  is a point in  $n$ -dimensional real space.

Alternatively, we can view  $\mathbf{x}_i$  as an  $n$ -dimensional vector (vectors in real space are considered to be columns):

$$\mathbf{x}_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{in} \end{pmatrix} = (x_{i1}, \dots, x_{in})^T.$$

**Linear Representation of vectors.** Let

$$I = \begin{pmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \dots & \mathbf{e}_n \\ 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix}$$

be a matrix of unit vectors

$$\mathbf{e}_i = (0, 0, \dots, 1_i, 0, \dots, 0)^T.$$

The vectors  $\mathbf{e}_1, \dots, \mathbf{e}_n$  are said to form the *standard basis* of  $\mathbb{R}^n$ .

Vector  $\mathbf{x}_i$  can be represented as a linear combination of  $e_1, \dots, e_n$  as

$$\mathbf{x}_i = x_{i1}\mathbf{e}_1 + x_{i2}\mathbf{e}_2 + \dots + x_{in}\mathbf{e}_n = \sum_{j=1}^n x_{ij}\mathbf{e}_j.$$

**More notation.** From the above, we can denote/represent the dataset  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  as follows:

$$X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{pmatrix} = \begin{pmatrix} -\mathbf{x}_1^T - \\ -\mathbf{x}_2^T - \\ \vdots \\ -\mathbf{x}_m^T - \end{pmatrix} = \left( \begin{array}{c|c|c|c} | & | & & | \\ A_1 & A_2 & \dots & A_n \\ | & | & & | \end{array} \right)$$

## Data Point Manipulations

Let  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subseteq \mathbb{R}^n$  be a dataset of numeric data points, and let  $\mathbf{x} \in X$  and  $\mathbf{y} \in X$  be two data points from  $X$ :

$$\mathbf{x} = (x_1 \dots, x_n)^T$$

$$\mathbf{y} = (y_1 \dots, y_n)^T$$

We recall a number of important operations on individual  $n$ -dimensional vectors, and pairs of such vectors. This operations *play an important role in machine learning methodology*.

**Dot product.** The dot product of  $\mathbf{x}$  and  $\mathbf{y}$ , denoted  $\mathbf{x} \cdot \mathbf{y}$  is defined as

$$\begin{aligned}\mathbf{x} \cdot \mathbf{y} &= \mathbf{x}^T \mathbf{y} = (x_1, \dots, x_n) \times \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \\ &= x_1 y_1 + x_2 y_2 + \dots + x_n y_n = \sum_{i=1}^n x_i y_i.\end{aligned}$$

**Orthogonality.** Two vectors  $\mathbf{x}$  and  $\mathbf{y}$  are called **orthogonal** iff

$$\mathbf{x} \cdot \mathbf{y} = 0.$$

**Vector norms.** The  $L_1$ -norm of a vector  $\mathbf{x}$  is defined as

$$\|\mathbf{x}\|_1 = |x_1| + |x_2| + \dots + |x_n| = \sum_{i=1}^n |x_i|.$$

The  $L_2$ -norm of vector  $\mathbf{x}$ , also called the *Euclidean norm* or the *length* of  $\mathbf{x}$  is defined as

$$\|\mathbf{x}\|_2 = \|\mathbf{x}\| = \sqrt{\mathbf{x}^T \cdot \mathbf{x}} = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2} = \sqrt{\sum_{i=1}^n x_i^2}.$$

In general, an  $L_p$ -norm of a vector  $\mathbf{x}$  is defined as

$$\|\mathbf{x}\|_p = (x_1^p + x_2^p + \dots + x_n^p)^{\frac{1}{p}} = \left( \sum_{i=1}^n x_i^p \right)^{\frac{1}{p}}.$$

**Distance.** The *Euclidean distance* between two vectors  $\mathbf{x}$  and  $\mathbf{y}$  is defined as

$$\delta(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\| = \sqrt{(\mathbf{x} - \mathbf{y})^T (\mathbf{x} - \mathbf{y})} = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}.$$

Similarly, the  $L_p$ -norm distance between  $\mathbf{x}$  and  $\mathbf{y}$  is defined as

$$\delta_p(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_p$$

In particular, the  $L_1$ -norm distance between  $\mathbf{x}$  and  $\mathbf{y}$ , also known as the *Manhattan distance* is

$$\text{Manhattan}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_1 = |x_1 - y_1| + |x_2 - y_2| + \dots + |x_n - y_n| = \sum_{i=1}^n |x_i - y_i|.$$



**Angle.** In many settings, it is helpful to compare the *directions* of the  $n$ -dimensional vectors  $\mathbf{x}$  and  $\mathbf{y}$  from the dataset  $X$ . In such cases, we use the **cosine similarity** score, which is the cosine of the angle between the two vectors. It is defined as follows:

$$\begin{aligned}\cos(\mathbf{x}, \mathbf{y}) &= \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} = \left( \frac{\mathbf{x}}{\|\mathbf{x}\|} \right)^T \left( \frac{\mathbf{y}}{\|\mathbf{y}\|} \right) = \\ &= \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}.\end{aligned}$$

**Cauchy-Schwartz inequality.** The *Cauchy-Schwartz inequality* for the  $L_2$ -norm and the dot product of two vectors state that

$$|\mathbf{x} \cdot \mathbf{y}| \leq \|\mathbf{x}\| \|\mathbf{y}\|$$

Therefore,

$$-1 \leq \cos(\mathbf{x}, \mathbf{y}) \leq 1$$

**Mean point of a dataset.** Given a dataset  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ , the mean point of the dataset or the mean is defined as:

$$\text{mean}(X) = \mu_{\mathbf{X}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$$

**Total Variance.** The total variance of the dataset  $X$  is *the average squared distance from a point in  $X$  to the mean point of  $X$* :

$$\text{var}(X) = \frac{1}{m} \sum_{i=1}^m \delta(\mathbf{x}_i, \mu_{\mathbf{X}}) = \frac{1}{m} \sum_{i=1}^m \|\mathbf{x}_i - \mu_{\mathbf{X}}\|^2$$

We can simplify the computation of variance as follows:

$$\begin{aligned}\text{var}(X) &= \frac{1}{m} \sum_{i=1}^m \|\mathbf{x}_i - \mu_{\mathbf{X}}\|^2 = \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n (x_{ij} - \mu_{Xj})^2 = \\ &= \frac{1}{m} \sum_{i=1}^m (\|\mathbf{x}_i\|^2 - 2\mathbf{x}_i^T \mu_{\mathbf{X}} + \|\mu_{\mathbf{X}}\|^2) = \frac{1}{m} \left( \sum_{i=1}^m \|\mathbf{x}_i\|^2 - 2m\mu_{\mathbf{X}}^T \left( \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i \right) + m\|\mu_{\mathbf{X}}\|^2 \right) = \\ &= \frac{1}{m} \left( \sum_{i=1}^m \|\mathbf{x}_i\|^2 - 2m\mu_{\mathbf{X}}^T \mu_{\mathbf{X}} + m\|\mu_{\mathbf{X}}\|^2 \right) = \\ &= \frac{1}{m} \left( \sum_{i=1}^m \|\mathbf{x}_i\|^2 \right) - \|\mu_{\mathbf{X}}\|^2\end{aligned}$$

**Orthogonal Projection.** Consider two vectors (data points)  $\mathbf{x}, \mathbf{y} \in X$ .

An **orthogonal decomposition** of  $\mathbf{x}$  in the direction of  $\mathbf{y}$  is a pair of vectors  $\mathbf{p}, \mathbf{r}$ , such that the following holds:

$$\begin{aligned}\mathbf{x} &= \mathbf{p} + \mathbf{r} \\ \mathbf{p}^T \mathbf{r} &= 0 \\ \mathbf{r}^T \mathbf{y} &= 0\end{aligned}$$

That is,  $\mathbf{x}$  can be decomposed as the sum of two *orthogonal vectors*,  $\mathbf{p}$  and  $\mathbf{r}$ , one of which ( $\mathbf{r}$ ) is, in turn, *orthogonal to*  $\mathbf{y}$  (which forces the other vector,  $\mathbf{p}$  to be co-aligned, or *colinear* with  $\mathbf{y}$ )

Vector  $\mathbf{p}$  is the **orthogonal projection** of  $\mathbf{x}$  onto  $\mathbf{y}$ . Vector  $\mathbf{r} = \mathbf{x} - \mathbf{p}$  is the **perpendicular distance** between  $\mathbf{x}$  and  $\mathbf{y}$ , and is colinear with the *normal vector* to  $\mathbf{y}$ .

Because  $\mathbf{p}$  is colinear with  $\mathbf{y}$ , we have

$$\mathbf{p} = c\mathbf{y},$$

and

$$\mathbf{r} = \mathbf{x} - c\mathbf{y}.$$

Because  $\mathbf{p}^T \mathbf{r} = 0$ , we get

$$\mathbf{p}^T \mathbf{r} = (c\mathbf{y})^T (\mathbf{x} - c\mathbf{y}) = c\mathbf{y}^T \mathbf{x} - c^2 \mathbf{y}^T \mathbf{y} = 0$$

Because  $c \neq 0$ , we must therefore have

$$\mathbf{y}^T \mathbf{x} - c\mathbf{y}^T \mathbf{y} = 0,$$

which yields

$$c = \frac{\mathbf{y}^T \mathbf{x}}{\mathbf{y}^T \mathbf{y}} = \frac{\mathbf{y} \cdot \mathbf{x}}{\|\mathbf{y}\|^2}$$

Therefore, the orthogonal projection of  $\mathbf{x}$  onto  $\mathbf{y}$  can be computed as

$$\mathbf{p} = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{y}\|^2} \mathbf{y}.$$

Another form of this expression is

$$\mathbf{p} = (\|\mathbf{x}\|^2 \cos(\mathbf{x}, \mathbf{y})) \mathbf{x}.$$

## More Linear Algebra Concepts

**Linear combinations.** Let  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$  be a dataset of points  $\mathbf{x}_i \in \mathbb{R}^n$ . Given a set of scalar values  $c_1, \dots, c_m$ , the vector

$$c_1\mathbf{x}_1 + c_2\mathbf{x}_2 + \dots + c_m\mathbf{x}_m$$

is called a *linear combination* of vectors  $\mathbf{x}_1, \dots, \mathbf{x}_m$ .

**Spanning set.** The set of all possible linear combinations

$$\mathbf{v} = \sum_{i=1}^m c_i \mathbf{x}_i$$

of vectors  $\mathbf{x}_1, \dots, \mathbf{x}_m$ , denoted  $\text{span}(\mathbf{x}_1, \dots, \mathbf{x}_m)$  or  $\text{span}(X)$  is called the **spanning set** of  $\mathbf{x}_1, \dots, \mathbf{x}_m$ .

**Row and Column Space.** Given the dataset  $X$  represented as an  $m \times n$  matrix:

$$X = \left( \begin{array}{c|cccc} & X_1 & X_2 & \dots & X_n \\ \hline \mathbf{x}_1 & x_{11} & x_{12} & \dots & x_{1n} \\ \mathbf{x}_2 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_m & x_{m1} & x_{m2} & \dots & x_{mn} \end{array} \right)$$

the **column span** of  $X$ , denoted  $\text{col}(X)$  is defined as the

$$\text{col}(X) = \text{span}(X_1, \dots, X_n),$$

and the **row span** of  $X$ , denoted  $\text{row}(X)$  is defined as the

$$\text{row}(X) = \text{span}(\mathbf{x}_1, \dots, \mathbf{x}_m)$$

We can see that  $\text{col}(X) \subseteq \mathbb{R}^m$  and  $\text{row}(X) \subseteq \mathbb{R}^n$ .

**Linear dependence.** A set of vectors  $\mathbf{v}_1, \dots, \mathbf{v}_k$  are *linearly dependent* iff at least one vector can be represented as a linear combination of others, i.e., if there are scalars  $c_1, \dots, c_k$  such that at least one  $c_i \neq 0$ , and

$$c_1\mathbf{v}_1 + \dots + c_k\mathbf{v}_k = 0$$

Vectors  $\mathbf{v}_1, \dots, \mathbf{v}_k$  are *linearly independent* otherwise.

**Rank.** Let  $S \subseteq \mathbb{R}^d$ . A *basis set* for  $S$  is a set of linearly independent vectors  $B = \{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  such that  $\text{span}(\mathbf{v}_1, \dots, \mathbf{v}_k) = S$ .

If all vectors in  $B$  are pairwise orthogonal, i.e.,  $(\forall i, j \in 1 \dots k, i \neq j)(\mathbf{v}_i \cdot \mathbf{v}_j = 0)$ ,  $B$  is called an *orthogonal basis* of  $S$ .

If for each vector  $v_i \in B$ ,  $\|v_i\| = 1$ , and  $B$  is an orthogonal basis, then,  $B$  is called an *orthonormal basis*.

The *standard basis* for  $\mathbb{R}^d$  is the basis consisting of vectors

$$\begin{aligned} \mathbf{e}_1 &= (1, 0, \dots, 0)^T \\ \mathbf{e}_2 &= (0, 1, \dots, 0)^T \\ &\dots \\ \mathbf{e}_d &= (0, 0, \dots, 1)^T \end{aligned}$$

**Theorem.** Every basis of a set  $S$  of vectors has the same number of vectors.

The number of vectors in a basis of  $S$  is called the *dimension of  $S$* , denoted  $\dim(S)$ .

**Theorem.** Given a matrix

$$A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix},$$

$$\dim(\text{col}(A)) = \dim(\text{row}(A)).$$

This value is called the *rank* of the matrix  $A$ , denoted  $\text{rank}(A)$ .