# Correlated features
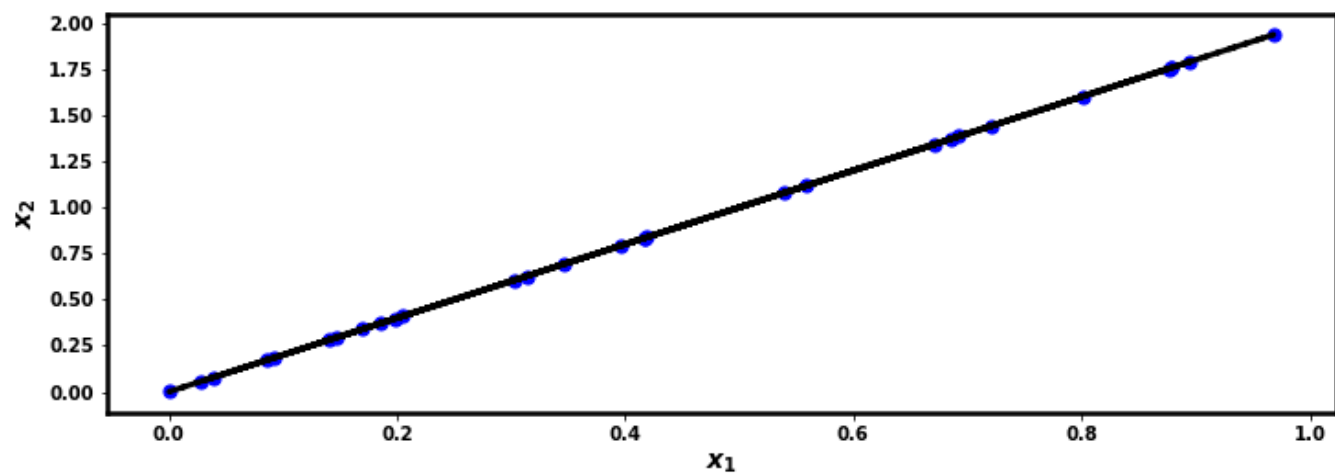
Consider the following set of examples with 2 features

As you can see

- $\backslash\mathbf{x}_2$ is perfectly correlated with $\backslash\mathbf{x}_1$
$$\backslash\mathbf{x}_2^{\backslash\mathrm{ip}} = 2 * \backslash\mathbf{x}_1^{\backslash\mathrm{ip}}$$

**Linear algebra**

A way to conceptualize $\x^{\ip}$

- As a point in the space spanned by unit basis vectors parallel to the horizontal and vertical axes.

$$\u_{(1)} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$\u_{(2)} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

- With $\x^{\ip}$ having exposure

$$\x_1^{\ip} \text{ to } \u_{(1)}$$

$$\x_2^{\ip} \text{ to } \u_{(2)}$$

So example $\x^{\ip}$ is

For example

$$\textcolor{red}{\backslash x}^{\backslash ip} \quad = \quad \begin{pmatrix} 3 \\ 6 \end{pmatrix}$$

$$= \quad 3 * \textcolor{red}{\backslash u}_{(1)} + 6 * \textcolor{red}{\backslash u}_{(2)}$$

$$= \quad 3 * \begin{pmatrix} 1 \\ 0 \end{pmatrix} + 6 * \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

$$= \quad \begin{pmatrix} 3 \\ 6 \end{pmatrix}$$

That is:

- Our feature space is defined by the basis vectors ("axes")

$$\backslash u_{(1)} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$\backslash u_{(2)} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

- $\backslash x^{\backslash ip}$ describes a point in the span of the basis vectors
  - $\backslash x_1^{\backslash ip}$ is the displacement of observation $\backslash x^{\backslash ip}$ along basis vector $\backslash u_{(1)}$
  - $\backslash x_2^{\backslash ip}$ is the displacement of observation $\backslash x^{\backslash ip}$ along basis vector $\backslash u_{(2)}$
- In general, for any length $n$ vector of features

$$\backslash x^{\backslash ip} = \sum_{j'=1}^{n} \backslash x_{j'}^{\backslash ip} * \backslash u_{(j')}$$

One could easily imagine a *different* set of basis vectors to describe the feature space

- For example: a rotation of basis vectors $\mathbf{u}_{(1)}, \ldots, \mathbf{u}_{(n)}$
- Let this alternate set of basis vectors be denoted by $\tilde{\mathbf{v}}_{(1)}, \ldots, \tilde{\mathbf{v}}_{(n)}$
- The basis vectors are mutually orthogonal

$$\tilde{\mathbf{v}}_{(1)} \cdot \tilde{\mathbf{v}}_{(2)} = 0$$

In the new basis space, example $\mathbf{x}^{\mathrm{ip}}$ has co-ordinates $\tilde{\mathbf{x}}^{\mathrm{ip}}$

$$\tilde{\mathbf{x}}^{\mathrm{ip}} = \sum_{j'=1}^{n} \tilde{\mathbf{x}}_{j'}^{\mathrm{ip}} * \tilde{\mathbf{v}}_{(j')}$$

PCA is a technique for finding particularly interesting alternate basis vectors.

The alternate basis is motivated by the fact that, for a given set of examples, there may be pairwise correlation among features.

- If the correlation is *perfect* for some pair of features, they are redundant
    - May drop one feature

Consider the set of examples above. Features 1 and 2 are perfectly correlated.

$$\backslash x_2^{\backslash ip} = 2 * \backslash x_1^{\backslash ip}$$

We can create an *alternate* basis vector (no longer parallel to the axes)

$$\backslash \tilde{v}_{(1)} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

such that example $\backslash x^{\backslash ip}$ has coordinates $\backslash \tilde{x}^{\backslash ip}$

$$\backslash \tilde{x}^{\backslash ip} = \backslash \tilde{x}_1^{\backslash ip} * \backslash \tilde{v}_{(1)}$$

Note that this alternate basis has only 1 basis vector, rather than the 2 basis vectors of the original representation.

For example

That is, $\x^{\ip}$ has exposure $\tilde{\x}_1^{\ip}$ to the new, single basis vector.

So

- Rather than representing $\x^{\ip}$ as a vector with 2 features
    - in the original basis
        - one basis vector per *raw* feature
        - mutually orthogonal basis vectors
- We can represent it as $\tilde{\x}^{\ip}$, a vector with 1 feature
    - in the new basis
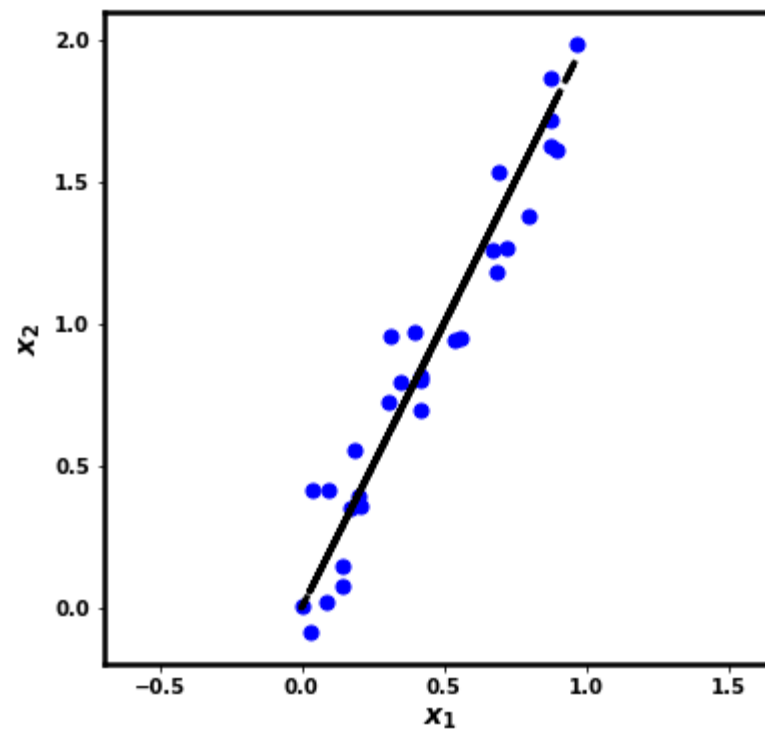    - which captures the correlation in 2 of the raw features when measured in the original basis

This is the essence of dimensionality reduction

- Changing bases to one with fewer basis vectors

It is rarely the case for features to be perfectly correlated

Let's modify the set of examples just a bit.

Two features: imperfect correlation

The single basis vector (black line)

- is insufficient to correctly capture each example
- *error*: displacement from black line

In order to eliminate the error, we add a *second* basis vector

- orthogonal to the first

So now $\backslash x^{\backslash ip}$ (measured in original basis) can be represented as $\tilde{\backslash x}^{\backslash ip}$ (measured in new basis)

- where
    - $\tilde{\backslash x}_{(1)}^{\backslash ip}$ measures the displacement along the first basis vector $\tilde{\backslash v}_1$
    - $\tilde{\backslash x}_{(2)}^{\backslash ip}$ measures the displacement along the second basis vector $\tilde{\backslash v}_2$

$$\tilde{\backslash x}^{\backslash ip} = \sum_{j'=1}^{2} \tilde{\backslash x}_{j'}^{\backslash ip} * \tilde{\backslash v}_{(j')}$$

- The dark black line in the diagram above is the first alternate basis vector $\tilde{\backslash v}_{(1)}$

In the diagram below, we add a second basis vector $\tilde{\backslash v}_{(2)}$

- orthogonal to the first

**Two features: imperfect correlation, alternate basis**

As you can see:

- The variation along $\tilde{\mathbf{v}}_{(1)}$ is much greater than that around $\tilde{\mathbf{v}}_{(2)}$
- Capturing the notion that the "main" relationship is along $\tilde{\mathbf{u}}_{(1)}$

In fact, if we dropped $\tilde{\mathbf{v}}_{(2)}$ such that $||\tilde{\mathbf{x}}|| = 1$

- The examples would be projected onto the line $\tilde{\mathbf{v}}_{(1)}$
- With little information being lost

PCA finds alternate basis vectors and *orders them* in order of decreasing variation.

# Subsets of correlated features

It may not be the case that a group of features is correlated across *all* examples

Consider our "equity factor model"

- consider two subsets of examples: stocks in/not in the "tech" sector
- all stocks in the first/second subset have the same loading on the "tech" factor (1/0)
- so there is correlation *within* the subsets but *not between* the subsets

Consider the MNIST digits

- The subset of examples corresponding to the digit "1"
- Have a particular set of correlated features (forming a vertical column of pixels near the middle of the image)
- Which *may not* be correlated with the same features in examples corresponding to *other* digits

Thus, a synthetic feature encodes a "concept" that occurs in many but not all examples

We will present a method to *discover* "concepts"

- It may not necessarily be the pattern of features that corresponds to an entire digit
- It may be a partial pattern common to several digits
    - Vertical band (0, 1, 4, 7)
    - Horizontal band at top (5, 7, 9)

```
In [5]: print("Done")
```

Done