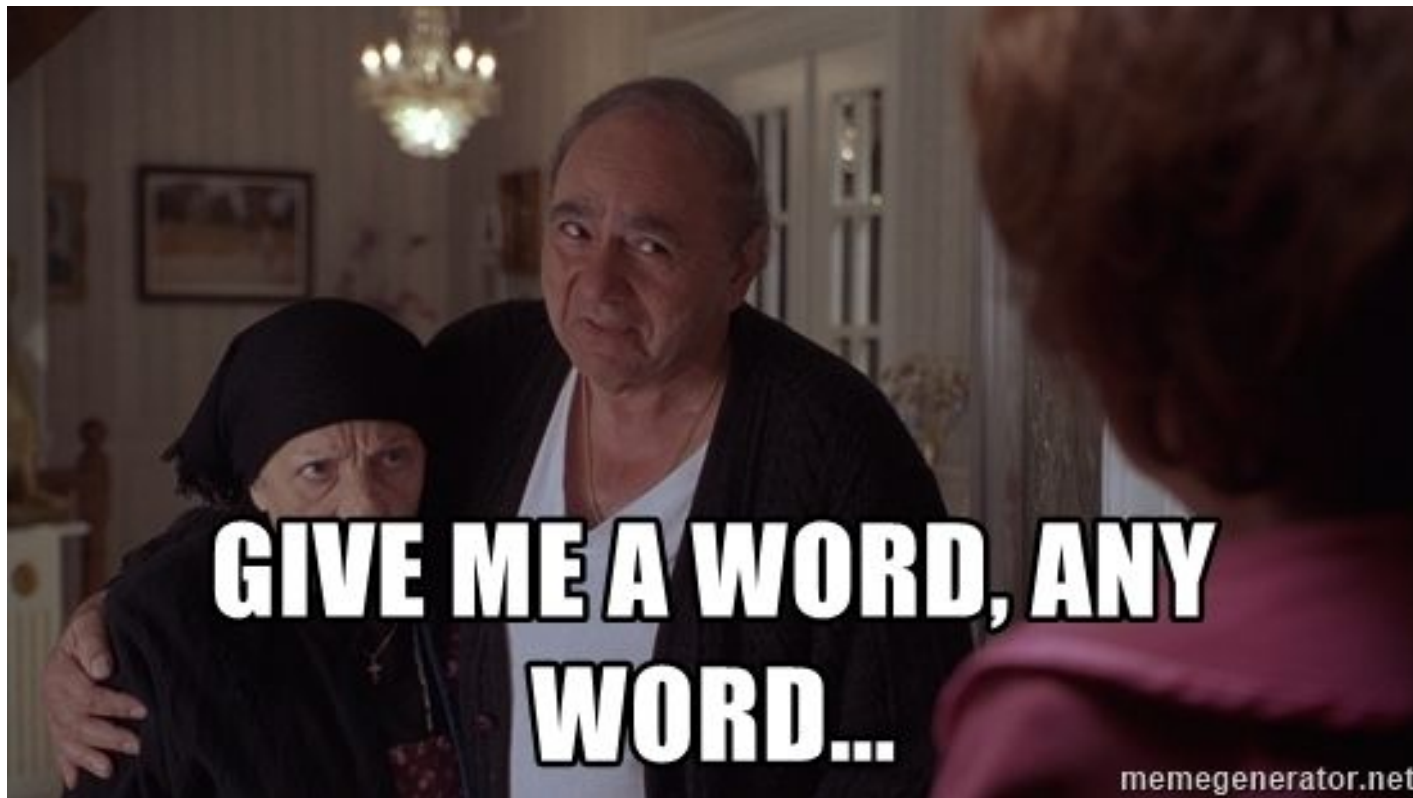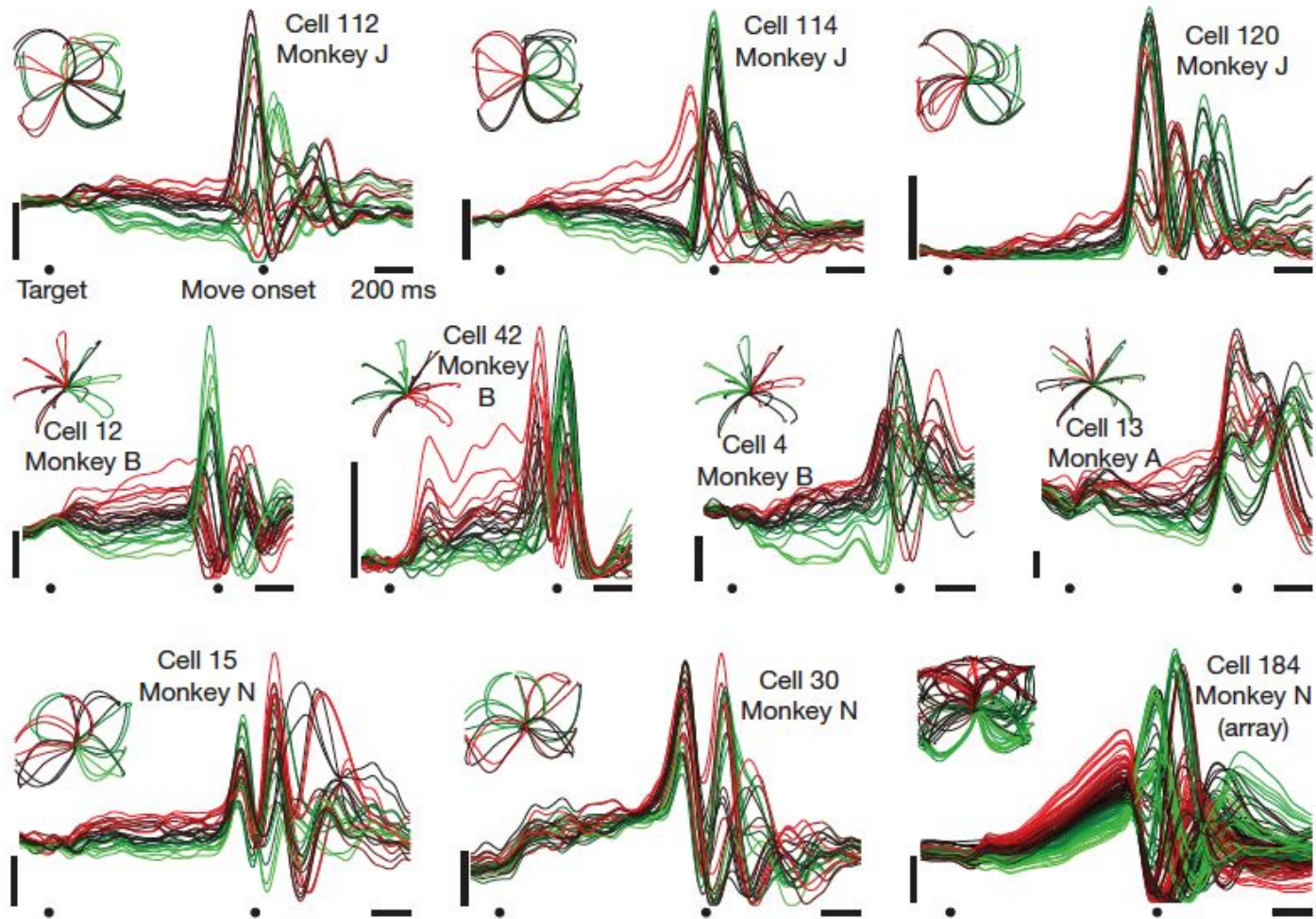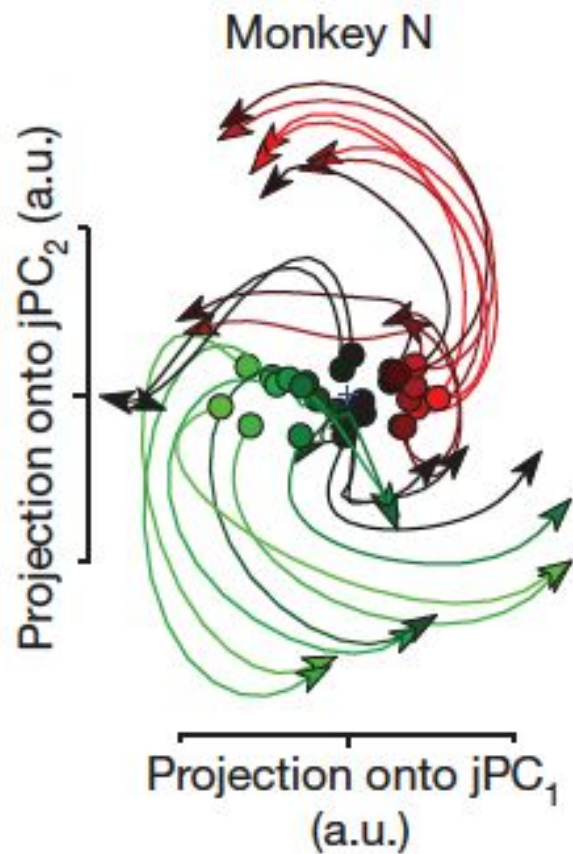# Principal component analysis

# Principal component analysis (PCA)

- workhorse of machine learning
- primarily used for dimensionality reduction
- unsupervised (does not predict)
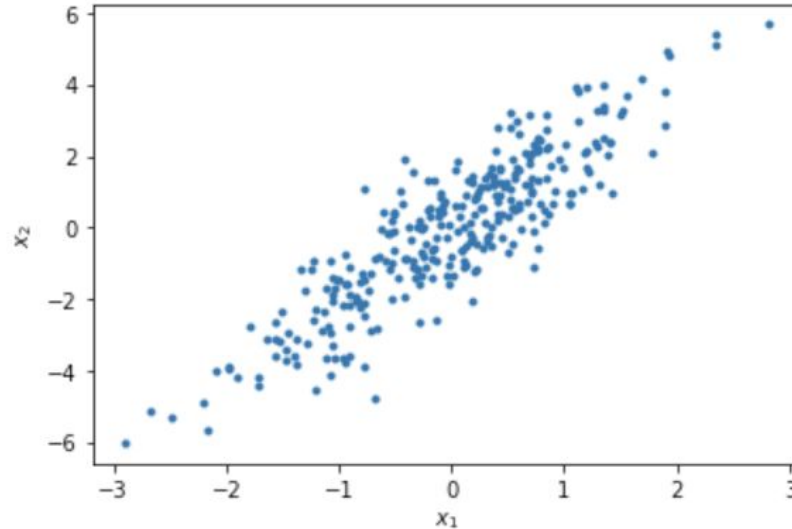- extension of "correlation" for multiple variables

Give me a machine learning algorithm,
and I can show you at the root of that algorithm is PCA!

Cell 112 Monkey J
Cell 114 Monkey J
Cell 120 Monkey J

Target     Move onset     200 ms

Cell 12 Monkey B
Cell 42 Monkey B
Cell 4 Monkey B
Cell 13 Monkey A

Cell 15 Monkey N
Cell 30 Monkey N
Cell 184 Monkey N (array)

Churchland et al., 2012

Monkey N          Monkey J-array          Monkey N-array

Projection onto jPC$_2$ (a.u.)

Projection onto jPC$_1$ (a.u.)

Churchland et al., 2012

Consider two variables, x1 and x2.

How can we "compress" them?

Consider two variables, x1 and x2.

How can we "compress" them?

Idea 1: Only keep x1 or x2 that has the largest variance.
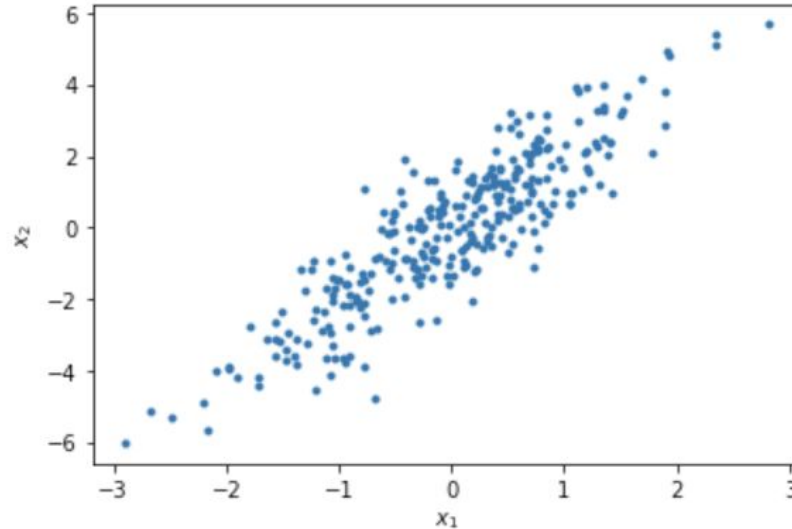
Consider two variables, x1 and x2.

How can we "compress" them?

Idea 1: Only keep x1 or x2 that has the largest variance.

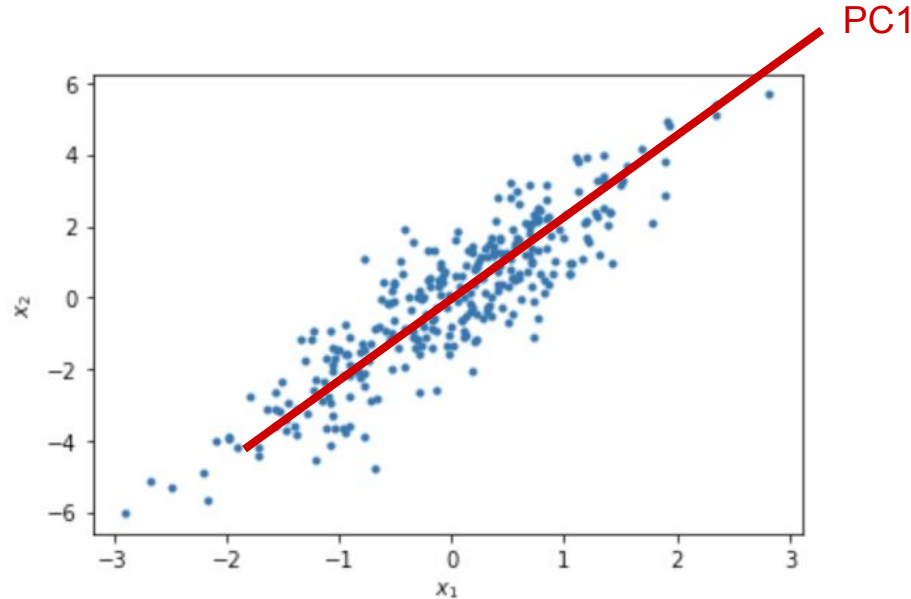Idea 2: Search for a combination of x1 and x2 that has the largest variance.

# Consider two variables, x1 and x2.

# How can we "compress" them?

# Consider two variables, x1 and x2.

# How can we "compress" them?

What happens for three variables: $x_1$, $x_2$, $x_3$?
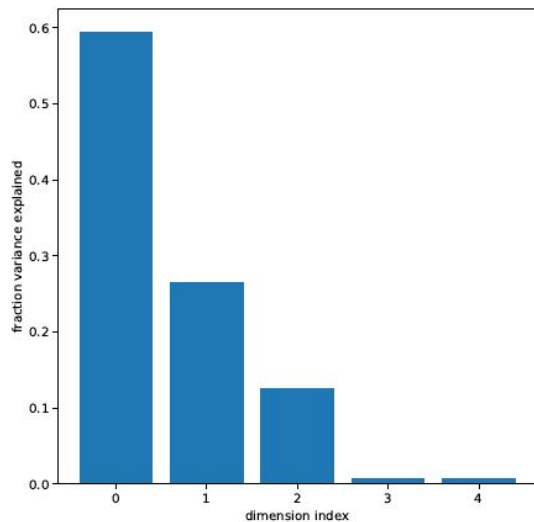
if data look like a pencil?

… a pancake?

… a sphere?
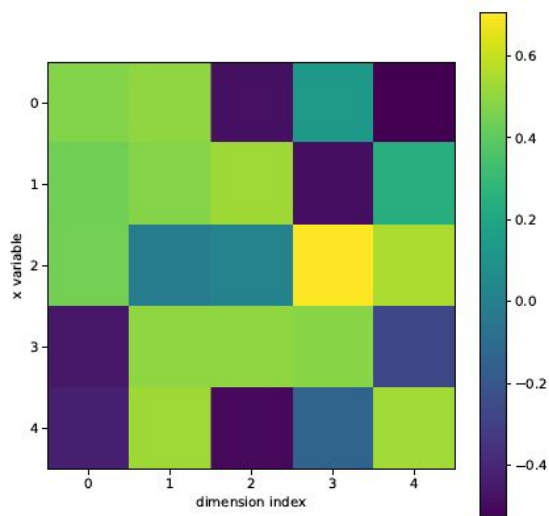
Three useful outputs of PCA:
- fraction variance explained
- loadings
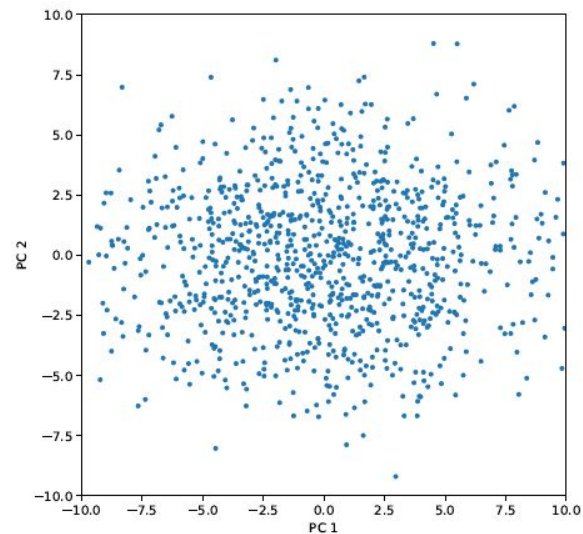- principal components

# Three useful outputs of PCA:

fraction explained variance

loadings

principal components

- Section 1 in Notebook

- Section 2 in Notebook