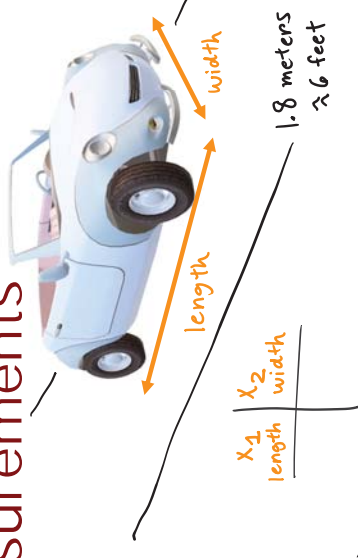# Principal Component Analysis (Optional)
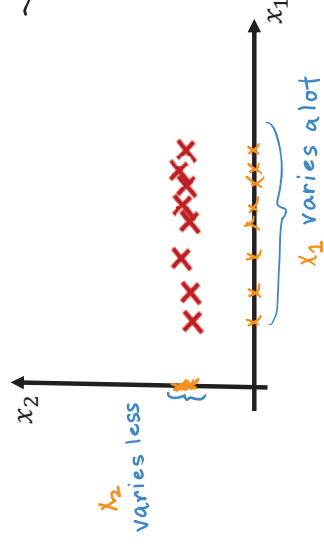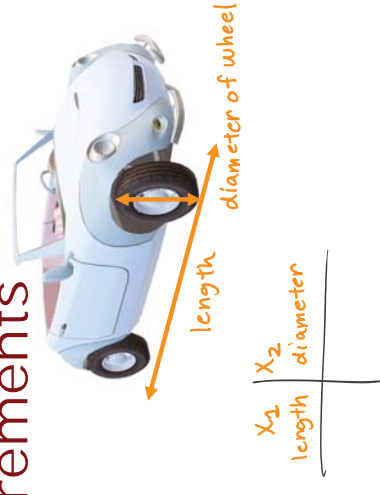
## Reducing the number of features

---

# Car measurements

width
length
1.8 meters ≈ 6 feet

$x_2$ width
$x_1$ length

can just take $x_1$ to reduce number of features

$x_2$

$x_1$

$x_1$ varies a lot

$x_2$ varies less

---

# Car measurements

height
length
width
diameter of wheel

---

# Car measurements

length
diameter of wheel

$x_2$ diameter
$x_1$ length

can just take $x_1$ to reduce number of features

$x_2$

$x_1$

$x_1$

$x_2$

## Size

PCA: find new axis and coordinates
use fewer numbers
to capture "size" feature

length

height

| | $x_1$ | $x_2$ |
| --- | --- | --- |
| | length | height |

2 → 1 features
many features → 2 or 3 features

z axis "size"

length

height

$x_1$

$x_2$

| Country | GDP (trillions of US$) |
| --- | --- |
| Canada | 1.577 |
| China | 5.878 |
| India | 1.632 |
| Russia | 1.48 |
| Singapore | 0.223 |
| USA | 14.527 |
| ... | ... |

$x_1$

---

| Country | GDP (trillions of US$) | Per capita GDP (thousands of intl. $) |
| --- | --- | --- |
| Canada | 1.577 | 39.17 |
| China | 5.878 | 7.54 |
| India | 1.632 | 3.41 |
| Russia | 1.48 | 19.84 |
| Singapore | 0.223 | 56.69 |
| USA | 14.527 | 46.86 |

$x_1$

$x_2$

---

## From 3D to 2D

3D

2D

$x_1$

$x_2$

$x_3$

$z_1$

$z_2$

| Country | GDP (trillions of US$) | Per capita GDP (thousands of intl. $) | Human Development Index |
|---|---|---|---|
| Canada | 1.577 | 39.17 | 0.908 |
| China | 5.878 | 7.54 | 0.687 |
| India | 1.632 | 3.41 | 0.547 |
| Russia | 1.48 | 19.84 | 0.755 |
| Singapore | 0.223 | 56.69 | 0.866 |
| USA | 14.527 | 46.86 | 0.91 |
| ... | ... | ... | ... |

$x_3$ $x_2$ $x_1$

50 features

| | $x_1$ GDP (trillions of US$) | $x_2$ Per capita GDP (thousands of intl. $) | $x_3$ Human Development Index | $x_4$ Life expectancy | ●●● Poverty Index (Gini as percentage) | Mean household income (thousands of US$) |
|---|---|---|---|---|---|---|
| Country | | | | | | |
| Canada | 1.577 | 39.17 | 0.908 | 80.7 | 32.6 | 67.293 |
| China | 5.878 | 7.54 | 0.687 | 73 | 46.9 | 10.22 |
| India | 1.632 | 3.41 | 0.547 | 64.7 | 36.8 | 0.735 |
| Russia | 1.48 | 19.84 | 0.755 | 65.5 | 39.9 | 0.72 |

| Country | $z_1$ | $z_2$ |
|---|---|---|
| Canada | 1.6 | 1.2 |
| China | 1.7 | 0.3 |
| India | 1.6 | 0.2 |
| Russia | 1.4 | 0.5 |

$z_2$ $z_1$

what if 50 features?

| Country | GDP (trillions of US$) | Per capita GDP (thousands of intl. $) | Human Development Index | Life expectancy |
|---|---|---|---|---|
| Canada | 1.577 | 39.17 | 0.908 | 80.7 |
| China | 5.878 | 7.54 | 0.687 | 73 |
| India | 1.632 | 3.41 | 0.547 | 64.7 |
| Russia | 1.48 | 19.84 | 0.755 | 65.5 |
| Singapore | 0.223 | 56.69 | 0.866 | 80 |
| USA | 14.527 | 46.86 | 0.91 | 78.3 |
| ... | ... | ... | ... | ... |

# Data visualization

← Singapore

↖ USA

50 features ≥ 2 features

reduce 50 "dimensional" data to 2D data

country's GDP

Per Person GDP

$z_2$

$z_2$

$z_1$

# Principal Component Analysis

## PCA Algorithm

---

## Choose an axis

z

origin

"project" examples onto the axis

Variance is large capturing info of original data

---

## Choose an axis

squished together less variance

not a good axis

z

---

## PCA algorithm

num bedrooms

$x_1$

size

$x_2$

10

8

coordinates $x_1 = 10$  $x_2 = 8$

Can we choose a different axis?

housing example: course1 week2

Preprocess features

Normalized to have zero mean

feature scaling

## More principal components

50 features $\rightarrow$ 3 principal components

$z_1$
$z_3$ — 3rd principal component
$z_2$
$90°$ $90°$

$z_1$
$z_2$ — 2nd principal component
at 90° angle "perpendicular"
1st principal component

---

## Choose an axis

principal component max variance

far apart capturing more info

---

## PCA is not linear regression

linear regression $(x, y)$
minimize distance along y axis
$y$
$x$

PCA
$x_1, x_2, \dots, x_{50}$ find axis to retain variance (info)
$z$
$x_2$
$x_1$

---

## Coordinate on the new axis

dot product

$$\begin{bmatrix} 2 \\ 3 \end{bmatrix} \cdot \begin{bmatrix} 0.71 \\ 0.71 \end{bmatrix}$$

$2 \times 0.71 + 3 \times 0.71 = 3.55$

$x_2$
coordinates $\begin{bmatrix} 2 \\ 3 \end{bmatrix}$
$z$
$3.55$
$x$
$2$
length 1 vector $\begin{bmatrix} 0.71 \\ 0.71 \end{bmatrix}$
$x_1$

## PCA is not linear regression

linear regression

$f(x) = 0$

PCA

$1^{st}$ principal component

---

## PCA in Code

PCA

DeepLearning.AI
Stanford ONLINE

---

## Approximation to the original data

original coordinates $\begin{bmatrix} 2 \\ 3 \end{bmatrix}$

2.52

3 —
2.52 —

z

3.55

given z = 3.55,
find original $(x_1, x_2)$ (approximately)
"reconstruction"

$$3.55 \times \begin{bmatrix} 0.71 \\ 0.71 \end{bmatrix} = \begin{bmatrix} 2.52 \\ 2.52 \end{bmatrix}$$

---

## PCA in scikit-learn

$z_3$ $z_1$ $z_2$

for visualization

info

Optional pre-processing: Perform feature scaling

1. "fit" the data to obtain 2 (or 3) new axes (principal components)

   fit includes mean normalization

2. Optionally examine how much variance is explained by each principal component.

   explained_variance_ratio

3. Transform (project) the data onto the new axes

   transform

# Applications of PCA

☆ Visualization *reduce to 2 or 3 features*

Less frequently used for:
- Data compression
  (to reduce storage or transmission costs) *50 → 10*
- Speeding up training of a supervised learning model

*n = 1000 → 100*

---

# Example

```
X = np.array([[1, 1], [2, 1], [3, 2],
              [-1, -1], [-2, -1], [-3, -2]])
```

2D

```
pca_1 = PCA(n_components=1)
pca_1.fit(X)
pca_1.explained_variance_ratio_   0.992
X_trans_1 = pca_1.transform(X)
X_reduced_1 = pca_1.inverse_transform(X_trans_1)
```

1D

```
array([
  [ 1.38340578],
  [ 2.22189802],
  [ 3.6053038 ],
  [-1.38340578],
  [-2.22189802],
  [-3.6053038 ]])
```

2D

---

# Example

```
X = np.array([[1, 1], [2, 1], [3, 2],
              [-1, -1], [-2, -1], [-3, -2]])
```

2D

```
pca_2 = PCA(n_components=2)
pca_2.fit(X)
pca_2.explained_variance_ratio_   0.992  0.008
X_trans_2 = pca_2.transform(X)
X_reduced_2 = pca_2.inverse_transform(X_trans_2)
```

$z_1$   $z_2$

```
array([
  [ 1.38340578,  0.2935787 ],
  [ 2.22189802, -0.251334841],
  [ 3.6053038 ,  0.042243851],
  [-1.38340578, -0.2935787 ],
  [-2.22189802,  0.251334841],
  [-3.6053038 , -0.042243851]])
```

2D