# 10-701 Introduction to Machine Learning

PCA

# PCA

# Raw data can be Complex, High-dimensional

To understand a phenomenon we measure various related quantities

If we knew what to measure or how to represent our measurements we might find simple relationships

But in practice we often *measure redundant signals*, e.g., US and European shoe sizes

We also *represent data via the method by which it was gathered*, e.g., pixel representation of brain imaging data

# Dimensionality Reduction

**Issues**
- *Measure redundant signals*
- *Represent data via the method by which it was gathered*

**Goal**: Find a 'better' representation for data
- To visualize and discover hidden patterns
- Preprocessing for supervised task

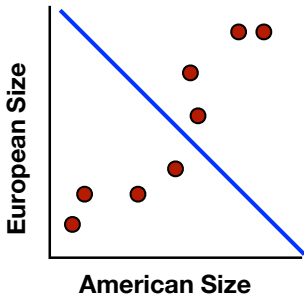How do we define 'better'?

# E.g., Shoe Size

We take noisy measurements on European and American scale
- Modulo noise, we expect perfect correlation

How can we do 'better', i.e., find a simpler, compact representation?
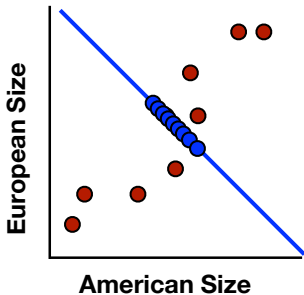- Pick a direction and project onto this direction

# E.g., Shoe Size

We take noisy measurements on European and American scale
- Modulo noise, we expect perfect correlation

How can we do 'better', i.e., find a simpler, compact representation?
- Pick a direction and project onto this direction

# E.g., Shoe Size

We take noisy measurements on European and American scale
- Modulo noise, we expect perfect correlation

How can we do 'better', i.e., find a simpler, compact representation?
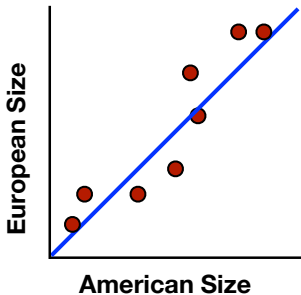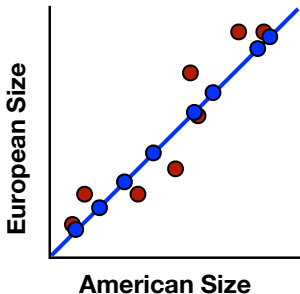- Pick a direction and project onto this direction

# E.g., Shoe Size

We take noisy measurements on European and American scale
- Modulo noise, we expect perfect correlation

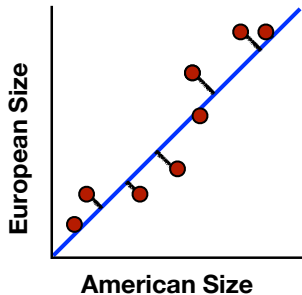How can we do 'better', i.e., find a simpler, compact representation?
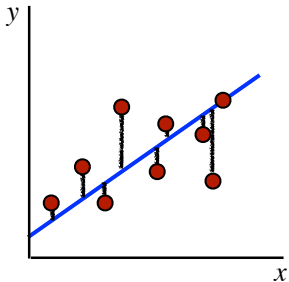- Pick a direction and project onto this direction

# Goal: Minimize Reconstruction Error

Minimize Euclidean distances between original points and their projections

PCA solution solves this problem!

**Linear Regression** — predict $y$ from $x$. Evaluate accuracy of predictions (represented by blue line) by **vertical** distances between points and the line

**PCA** — reconstruct 2D data via 2D data with single degree of freedom. Evaluate reconstructions (represented by blue line) by **Euclidean** distances

# Another Goal: Maximize Variance

To identify patterns we want to study variation across observations

Can we do 'better', i.e., find a compact representation that captures variation?

# Another Goal: Maximize Variance

To identify patterns we want to study variation across observations

Can we do 'better', i.e., find a compact representation that captures variation?

# Another Goal: Maximize Variance

To identify patterns we want to study variation across observations

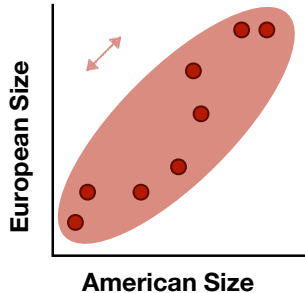Can we do 'better', i.e., find a compact representation that captures variation?

PCA solution finds directions of maximal variance!

# PCA Formulation

PCA: find lower-dimensional representation of raw data

- $\mathbf{X}$ is $n \times d$ (raw data)
- $\mathbf{Z} = \mathbf{XP}$ is $n \times k$ (reduced representation, PCA 'scores')
- $\mathbf{P}$ is $d \times k$ (columns are $k$ principal components)
- Variance constraints

Linearity assumption ( $\mathbf{Z} = \mathbf{XP}$ ) simplifies problem

$$\begin{bmatrix} \mathbf{Z} \end{bmatrix} = \begin{bmatrix} \mathbf{X} \end{bmatrix} \begin{bmatrix} \mathbf{P} \end{bmatrix}$$

Given $n$ training points with $d$ features:

- $\mathbf{X} \in \mathbb{R}^{n \times d}$: matrix storing points
- $x_j^{(i)}$: $j$th feature for $i$th point
- $\mu_j$ : mean of $j$th feature

Variance of 1st feature $\qquad \sigma_1^2 = \dfrac{1}{n} \displaystyle\sum_{i=1}^{n} \left( x_1^{(i)} - \mu_1 \right)^2$

Variance of 1st feature (assuming zero mean) $\qquad \sigma_1^2 = \dfrac{1}{n} \displaystyle\sum_{i=1}^{n} \left( x_1^{(i)} \right)^2$

Given $n$ training points with $d$ features:

- $\mathbf{X} \in \mathbb{R}^{n \times d}$ : matrix storing points
- $x_j^{(i)}$ : $j$th feature for $i$th point
- $\mu_j$ : mean of $j$th feature

Covariance of 1st and 2nd features (assuming zero mean)

$$\sigma_{12} = \frac{1}{n} \sum_{i=1}^{n} x_1^{(i)} x_2^{(i)}$$

- Symmetric: $\sigma_{12} = \sigma_{21}$
- Zero → uncorrelated
- Large magnitude → (anti) correlated / redundant
- $\sigma_{12} = \sigma_1^2 = \sigma_2^2$ → features are the same

# Covariance Matrix

Covariance matrix generalizes this idea for many features

$d \times d$ covariance matrix with zero mean features
$$\mathbf{C_X} = \frac{1}{n}\mathbf{X}^\top\mathbf{X}$$

- $i$th diagonal entry equals variance of $i$th feature
- $ij$th entry is covariance between $i$th and $j$th features
- Symmetric (makes sense given definition of covariance)

# PCA Formulation

PCA: find lower-dimensional representation of raw data
- $\mathbf{X}$ is $n \times d$ (raw data)
- $\mathbf{Z} = \mathbf{XP}$ is $n \times k$ (reduced representation, PCA 'scores')
- $\mathbf{P}$ is $d \times k$ (columns are $k$ principal components)
- Variance / Covariance constraints

What constraints make sense in reduced representation?
- No feature correlation, i.e., all off-diagonals in $\mathbf{C_Z}$ are zero
- Rank-ordered features by variance, i.e., sorted diagonals of $\mathbf{C_Z}$

# PCA Formulation

PCA: find lower-dimensional representation of raw data

- $\mathbf{X}$ is $n \times d$ (raw data)
- $\mathbf{Z} = \mathbf{XP}$ is $n \times k$ (reduced representation, PCA 'scores')
- $\mathbf{P}$ is $d \times k$ (columns are $k$ principal components)
- Variance / Covariance constraints

$\mathbf{P}$ equals the top $k$ eigenvectors of $\mathbf{C_X}$

$$\mathbf{Z} = \mathbf{X} \, \mathbf{P}$$

# PCA Solution

All covariance matrices have an eigendecomposition

- $\mathbf{C_X} = \mathbf{U\Lambda U}^\top$ (eigendecomposition)
- $\mathbf{U}$ is $d \times d$ (column are eigenvectors, sorted by their eigenvalues)
- $\mathbf{\Lambda}$ is $d \times d$ (diagonals are eigenvalues, off-diagonals are zero)

The $d$ eigenvectors are orthonormal directions of max variance

- Associated eigenvalues equal variance in these directions
- 1st eigenvector is direction of max variance (variance is $\lambda_1$)

# Choosing $k$

How should we pick the dimension of the new representation?

**Visualization**: Pick top 2 or 3 dimensions for plotting purposes

**Other analyses**: Capture 'most' of the variance in the data
- Recall that eigenvalues are variances in the directions specified by eigenvectors, and that eigenvalues are sorted

- Fraction of retained variance: $\dfrac{\sum_{i=1}^{k} \lambda_i}{\sum_{i=1}^{d} \lambda_i}$

Can choose $k$ such that we retain some fraction of the variance, e.g., 95%

# Other Practical Tips

PCA assumptions (linearity, orthogonality) not always appropriate
- Various extensions to PCA with different underlying assumptions, e.g., manifold learning, Kernel PCA, ICA

Centering is crucial, i.e., we must preprocess data so that all features have zero mean before applying PCA

PCA results dependent on scaling of data
- Data is sometimes rescaled in practice before applying PCA

# Orthogonal and Orthonormal Vectors

*Orthogonal* vectors are **perpendicular** to each other
- Equivalently, their dot product equals zero
- $\mathbf{a}^\top \mathbf{b} = 0$ and $\mathbf{d}^\top \mathbf{b} = 0$, but **c** isn't orthogonal to others



$$\mathbf{a} = \begin{bmatrix} 1 & 0 \end{bmatrix}^\top \qquad \mathbf{b} = \begin{bmatrix} 0 & 1 \end{bmatrix}^\top \qquad \mathbf{c} = \begin{bmatrix} 1 & 1 \end{bmatrix}^\top \qquad \mathbf{d} = \begin{bmatrix} 2 & 0 \end{bmatrix}^\top$$

*Orthonormal* vectors are orthogonal and have unit norm
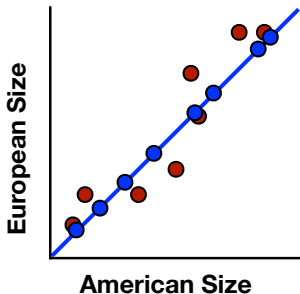- **a** are **b** are orthonormal, but **b** are **d** are not orthonormal

# PCA Iterative Algorithm

$k = 1$: Find direction of max variance, project onto this direction
● Locations along this direction are the new 1D representation

More generally, for $i$ in $\{1, \ldots, k\}$:
● Find direction of max variance that is *orthonormal* to previously selected directions, project onto this direction
● Locations along this direction are the $i$th feature in new representation

# Eigendecomposition

All covariance matrices have an eigendecomposition

- $\mathbf{C_X} = \mathbf{U \Lambda U}^\top$ (eigendecomposition)
- $\mathbf{U}$ is $d \times d$ (column are eigenvectors, sorted by their eigenvalues)
- $\mathbf{\Lambda}$ is $d \times d$ (diagonals are eigenvalues, off-diagonals are zero)

Eigenvector / Eigenvalue equation: $\mathbf{C_x u} = \lambda \mathbf{u}$

- By definition $\mathbf{u}^\top \mathbf{u} = 1$ (unit norm)

- Example: $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \implies$ eigenvector: $\mathbf{u} = \begin{bmatrix} 1 & 0 \end{bmatrix}^\top$
  eigenvalue: $\lambda = 1$

# PCA Formulation

PCA: find lower-dimensional representation of raw data

- $\mathbf{X}$ is $n \times d$ (raw data)
- $\mathbf{Z} = \mathbf{XP}$ is $n \times k$ (reduced representation, PCA 'scores')
- $\mathbf{P}$ is $d \times k$ (columns are $k$ principal components)
- Variance / Covariance constraints

$$\begin{bmatrix} \\ \mathbf{Z} \\ \\ \end{bmatrix} = \begin{bmatrix} \\ \mathbf{X} \\ \\ \end{bmatrix} \begin{bmatrix} \mathbf{P} \\ \\ \end{bmatrix}$$

# PCA Formulation, $k = 1$

PCA: find one-dimensional representation of raw data

- $\mathbf{X}$ is $n \times d$ (raw data)
- $\mathbf{z} = \mathbf{X}\mathbf{p}$ is $n \times 1$ (reduced representation, PCA 'scores')
- $\mathbf{p}$ is $d \times 1$ (columns are $k$ principal components)
- Variance constraint

$$\sigma_{\mathbf{z}}^2 = \frac{1}{n} \sum_{i=1}^{n} \left( z^{(i)} \right)^2 = ||\mathbf{z}||_2^2$$

**Goal**: Maximizes variance, i.e., $\max_{\mathbf{p}} \sigma_{\mathbf{z}}^2$ where $||\mathbf{p}||_2 = 1$

**Goal**: Maximizes variance, i.e., $\max_{\mathbf{p}} \sigma_{\mathbf{z}}^2$ where $||\mathbf{p}||_2 = 1$

$$\sigma_{\mathbf{z}}^2 = \frac{1}{n}||\mathbf{z}||_2^2$$

Relationship between Euclidean distance and dot product
$$= \frac{1}{n}\mathbf{z}^\top \mathbf{z}$$

Definition: $\mathbf{z} = \mathbf{Xp}$
$$= \frac{1}{n}(\mathbf{Xp})^\top(\mathbf{Xp})$$

Transpose property: $(\mathbf{Xp})^\top = \mathbf{p}^\top \mathbf{X}^\top$; associativity of multiply
$$= \frac{1}{n}\mathbf{p}^\top \mathbf{X}^\top \mathbf{Xp}$$

Definition: $\mathbf{C_X} = \frac{1}{n}\mathbf{X}^\top \mathbf{X}$
$$= \mathbf{p}^\top \mathbf{C_X} \mathbf{p}$$

**Restated Goal:** $\max_{\mathbf{p}} \mathbf{p}^\top \mathbf{C_x} \mathbf{p}$ where $||\mathbf{p}||_2 = 1$

# Connection to Eigenvectors

Recall eigenvector / eigenvalue equation: $\mathbf{C_x u} = \lambda \mathbf{u}$

- By definition $\mathbf{u}^\top \mathbf{u} = 1$, and thus $\mathbf{u}^\top \mathbf{C_x u} = \lambda$
- But this is the expression we're optimizing, and thus maximal variance achieved when $\mathbf{p}$ is top eigenvector of $\mathbf{C_X}$

Similar arguments can be used for $k > 1$

**Restated Goal:** $\max_{\mathbf{p}} \mathbf{p}^\top \mathbf{C_x p}$ where $||\mathbf{p}||_2 = 1$