# PROPOSAL GRADES POSTED

**Review Grades for Project Proposal**

● Regrade Requests Open    ● Grades Published



Histogram of score distribution

| Minimum | Median | Maximum | Mean | Std Dev |
|---------|--------|---------|------|---------|
| **95.0** | **98.5** | **100.0** | **98.0** | **1.51** |

# SPOTLIGHT REMINDER
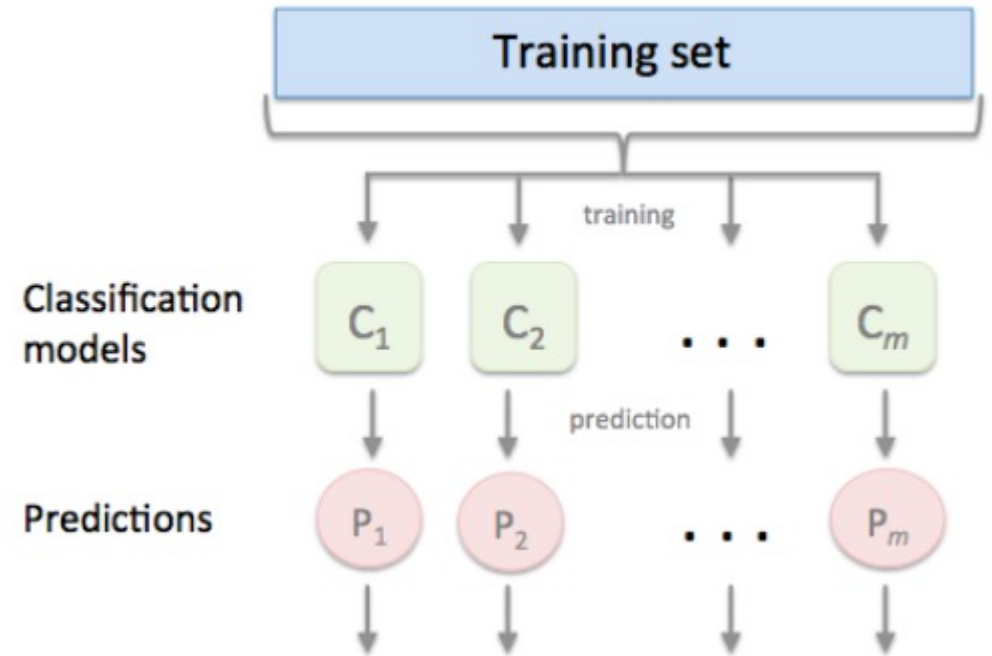
- Spotlight PPT slide due 10/30 on Canvas

- Project spotlight presentation on 11/1

  - 2 min each group

  - 1 min between groups

  - Group order 16 -> 1

# ENSEMBLE METHODS

- Ensemble of same classifiers

  - Bagging and random forest

  - Boosting and gradient boosted tree
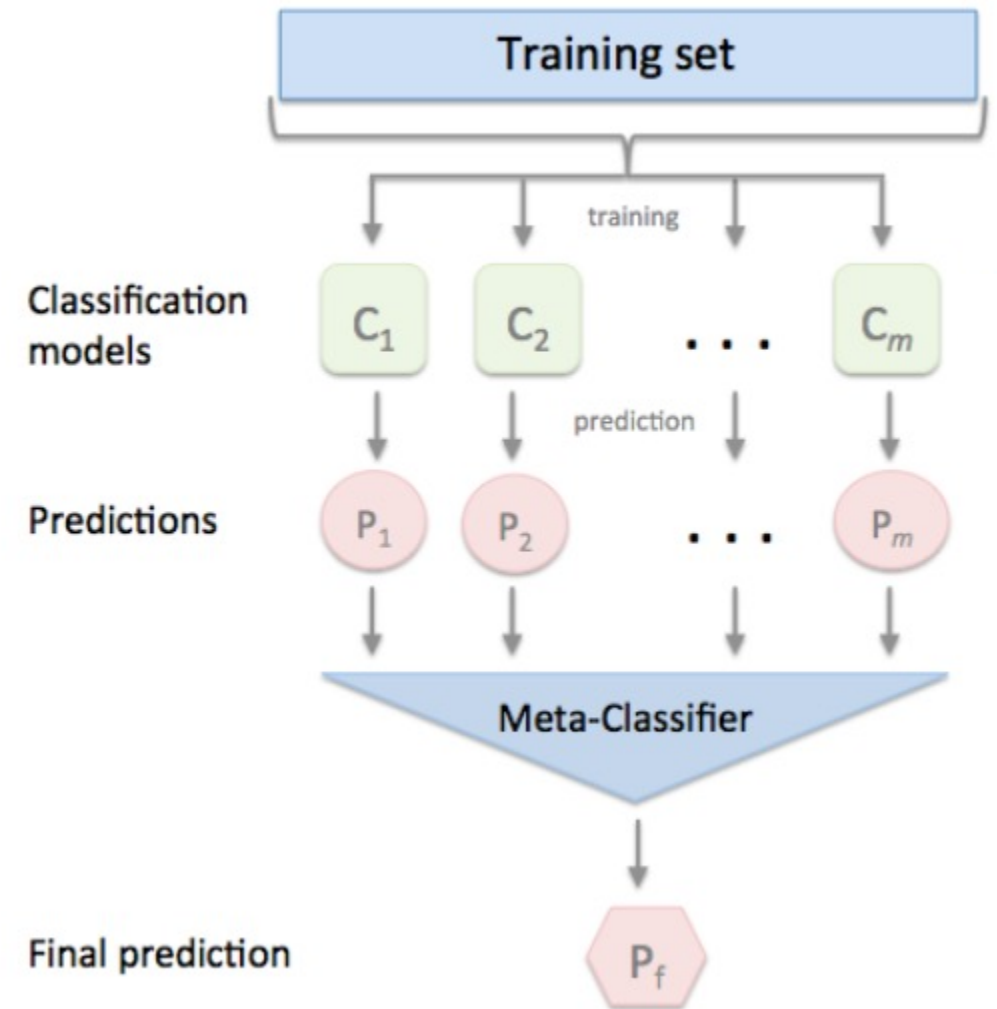
- Ensemble of different classifiers

# VOTING

- Voting ensembles mimic error-correcting codes

  - More voters —> potential better signal to noise

  - Lower correlation between models



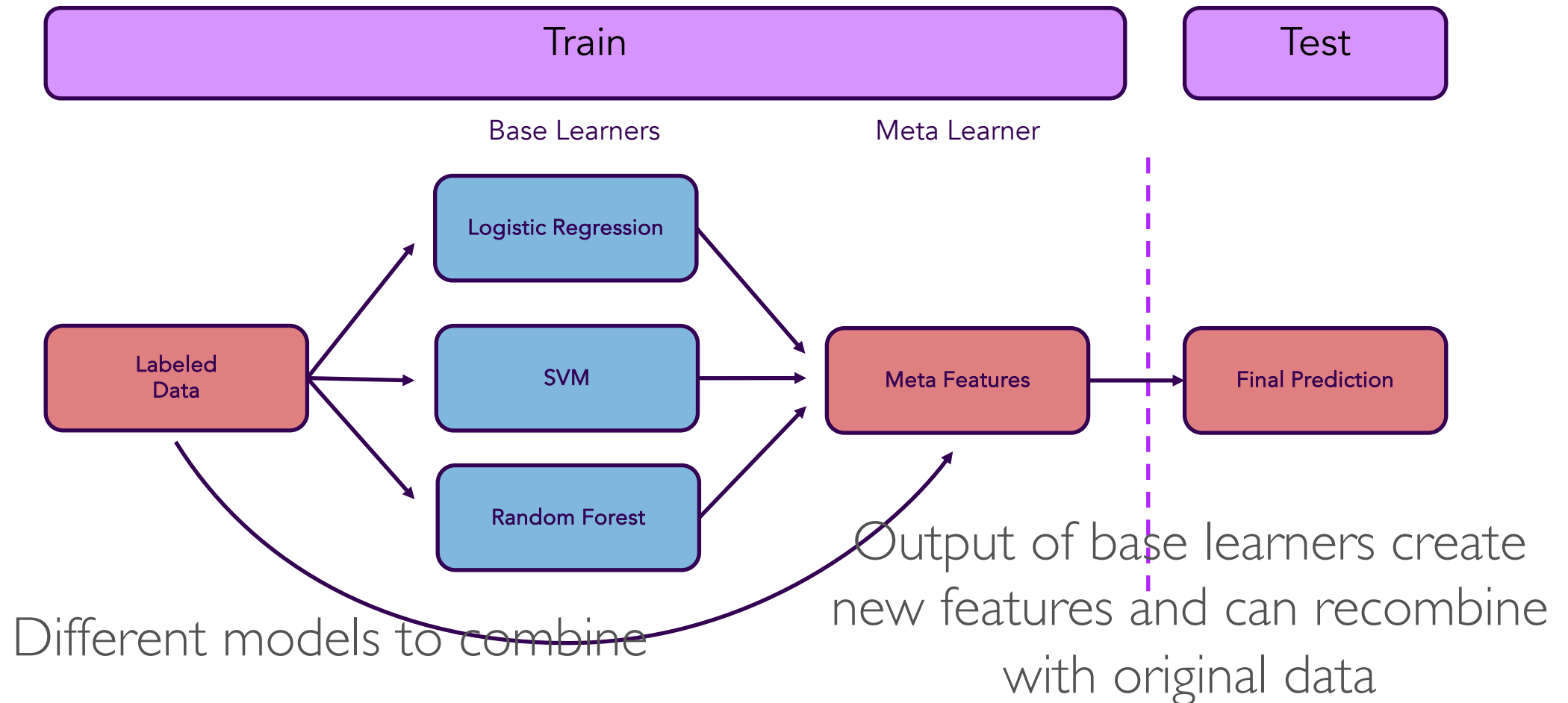How do we combine the predictions?

# STACKING

- Introduced by Wolpert, 1992

- Use a pool of base classifiers and do <span style="color:red">out-of-fold</span> predictions, then train a meta-classifier to combine their predictions

- Stacker model gains information by using first-stage predictions as features

# BLENDING

- Close to stacked generalization, but a bit simpler

- Instead of out-of-fold predictions, create small holdout set that the stacker is then trained on this set

- Disadvantages:

  - Less data used overall

  - Final model may overfit to holdout set
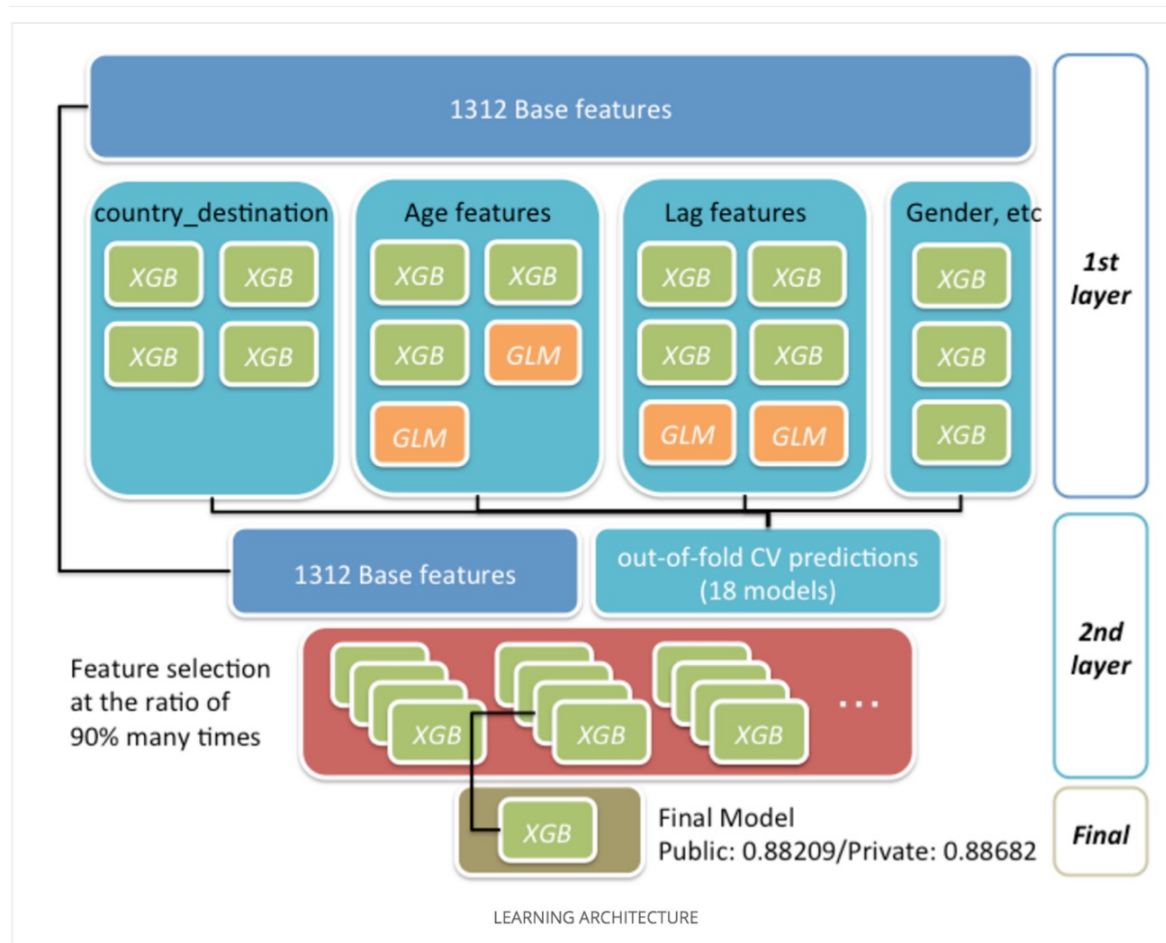
# EXAMPLE: COMBINING DIFFERENT CLASSIFIERS

# STACKING AND BLENDING

- Why stop at two stages? Why not combine multiple ensembles models?

# EXAMPLE: AIRBNB 2ND PLACE



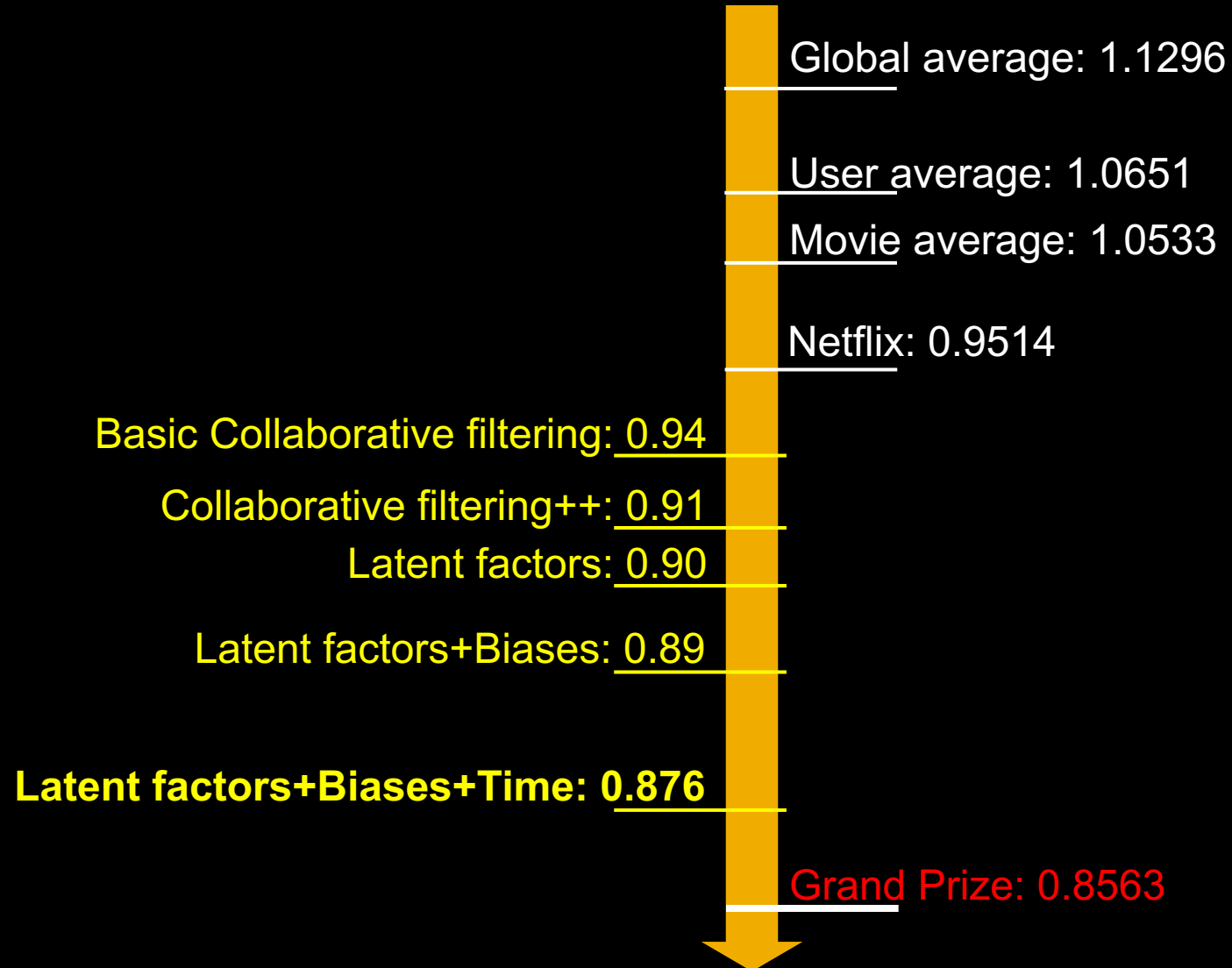https://github.com/Keiku/kaggle-airbnb-recruiting-new-user-bookings
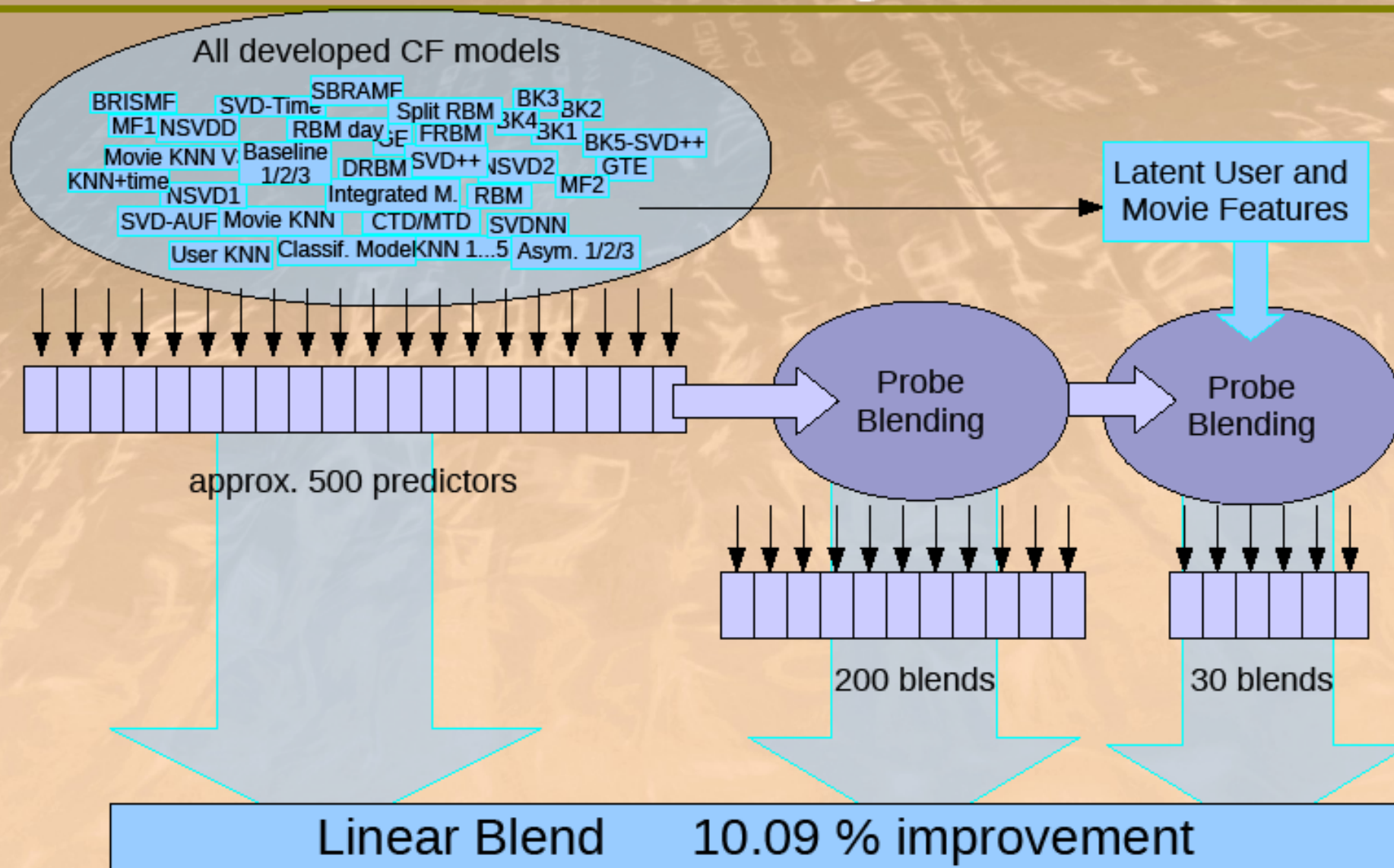
# The Netflix Prize

- An open competition 2006 – 2009 held by Netflix
- **Training data**
  - 100 million ratings (with timestamps), 480,000 users, 17,770 movies
  - 6 years of data: 2000-2005
- **Test data**
  - Last few ratings of each user (2.8 million)
- **Evaluation criterion:** Root Mean Square Error

(RMSE) $= \frac{1}{|R|} \sqrt{\sum_{(i,x) \in R} (\hat{r}_{xi} - r_{xi})^2}$

- **Competition**
  - **Netflix's system Cinematch RMSE: 0.9514**
  - **$1 million** prize for 10% improvement on Netflix
  - $50,000 progress prize every year

# Netflix Prize



Global average: 1.1296

User average: 1.0651

Movie average: 1.0533

Netflix: 0.9514

Basic Collaborative filtering: 0.94

Collaborative filtering++: 0.91

Latent factors: 0.90

Latent factors+Biases: 0.89

**Latent factors+Biases+Time: 0.876**

Grand Prize: 0.8563

# The big picture

## Solution of BellKor's Pragmatic Chaos

All developed CF models

BRISMF · SBRAMF · SVD-Time · Split RBM · BK3 · BK2 · MF1 · NSVDD · RBM day · GE · FRBM · BK4 · BK1 · BK5-SVD++ · Movie KNN v. · Baseline 1/2/3 · DRBM · SVD++ · NSVD2 · GTE · KNN+time · NSVD1 · Integrated M. · RBM · MF2 · SVD-AUF · Movie KNN · CTD/MTD · SVDNN · User KNN · Classif. Mode · KNN 1...5 · Asym. 1/2/3

Latent User and Movie Features

approx. 500 predictors

Probe Blending

Probe Blending

200 blends

30 blends

Linear Blend    10.09 % improvement

# DIMENSIONALITY REDUCTION

CS 334: Machine Learning

# TYPES OF UNSUPERVISED LEARNING

**Clustering**

identify unknown structure in the data

**Dimensionality Reduction**

use structural characteristics to simplify data

# CURSE OF DIMENSIONALITY

- Increasing features should improve performance right?

- In practice, too many features leads to worse performance

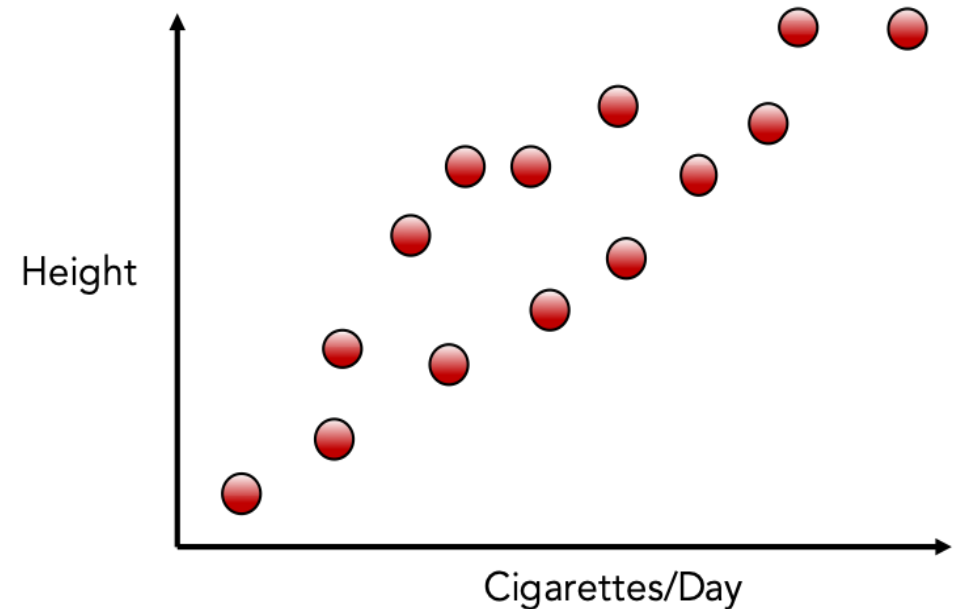- # of training examples required increases exponentially with the # of features



1 dimension: 10 positions

2 dimensions: 100 positions

3 dimensions: 1000 positions

# FEATURE SELECTION/REDUCTION

- Feature selection

  - Filter methods (agnostic to the learning algorithm)

  - Wrapper methods (keep model in the loop)

  - Embedded methods (e.g. Lasso regularization)
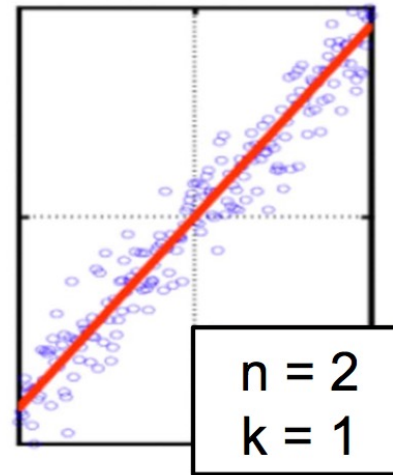
- Dimension reduction

# MOTIVATION: EXAMPLE

- Two features: Height and cigarettes per day

- Both features are correlated

- Can we reduce the features to one?

# DIMENSIONALITY REDUCTION

- Represent data with fewer dimensions

- Discover "intrinsic dimensionality" of data

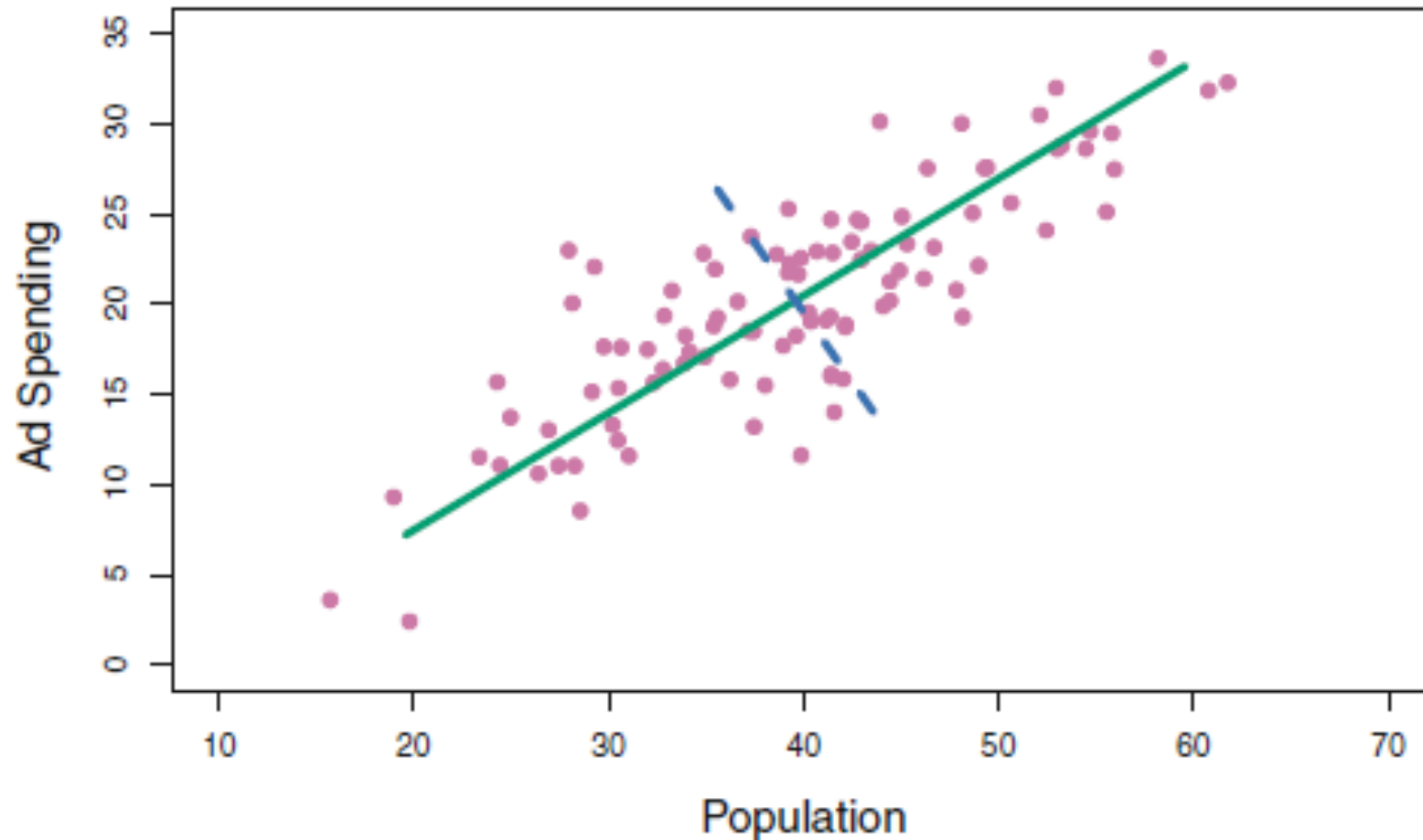- New feature space: linear / non-linear combination of original features

n = 2
k = 1

n = 3
k = 2

**Slide by Yi Zhang**

# WHY DIMENSIONALITY REDUCTION

- Simplified data processing

- Noise reduction (removing feature redundancy)

- Easier learning — less parameters

- Robust learning — numerical stability due to less correlations

- Easier visualization — show high dimensional data in 2D or 3D

GROUP ACTIVITY
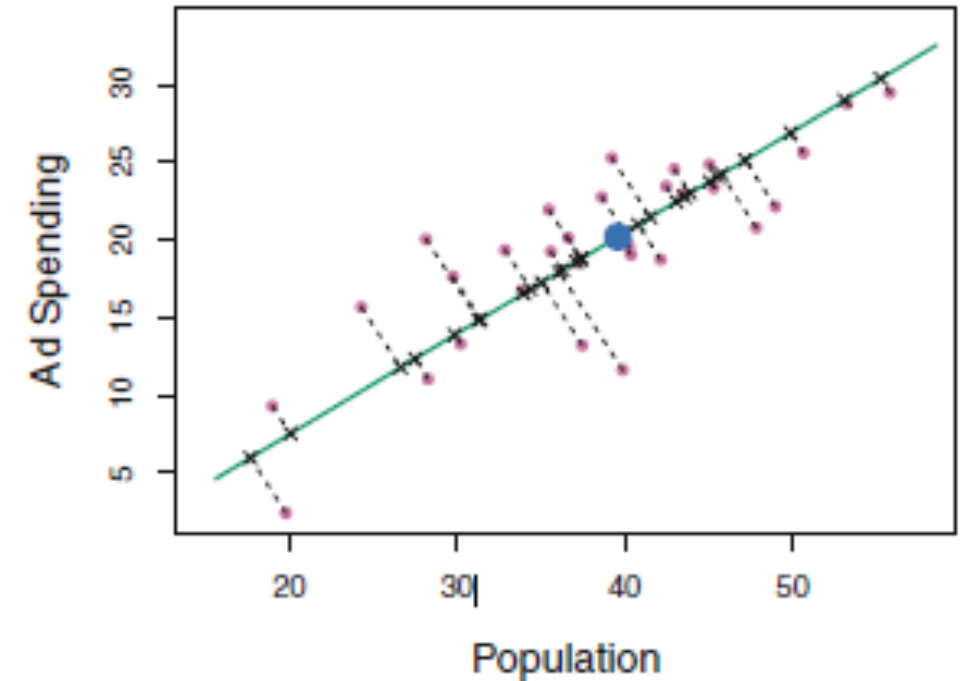
# WHICH IS THE BEST PROJECTION AND WHY?

# PRINCIPAL COMPONENT ANALYSIS (PCA)

- Developed by Pearson in 1901

- Popular and widely studied

- Finds sequence of linear combinations of the features (also known as principal components) that have maximal variance and are uncorrelated
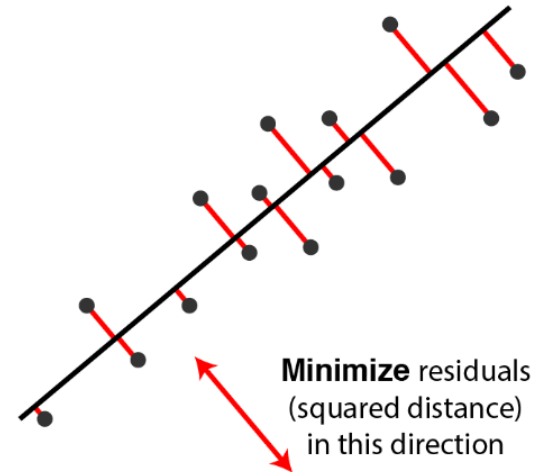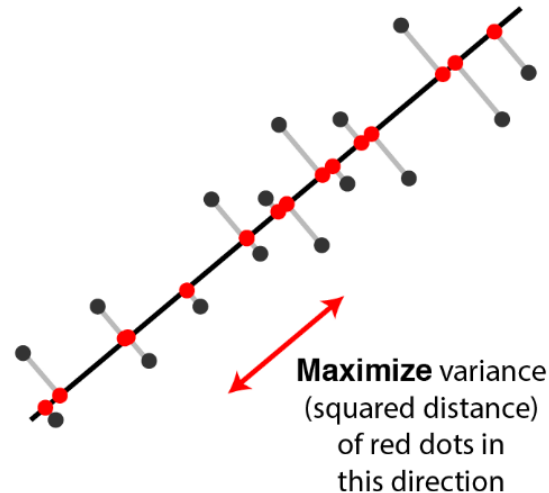
# PRINCIPAL COMPONENT ANALYSIS (PCA)

- First principal component

  - Yields the highest variance of the projection

  - Minimizes sum of squared perpendicular distances (reconstruction error between data and projected points)
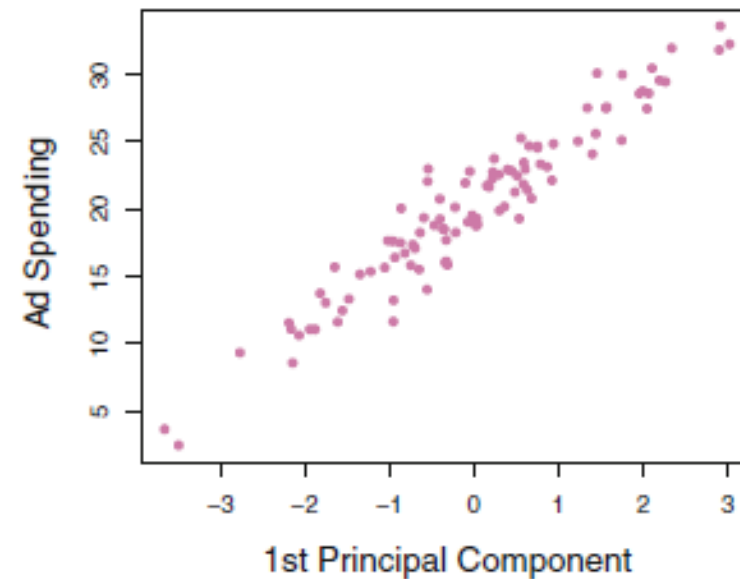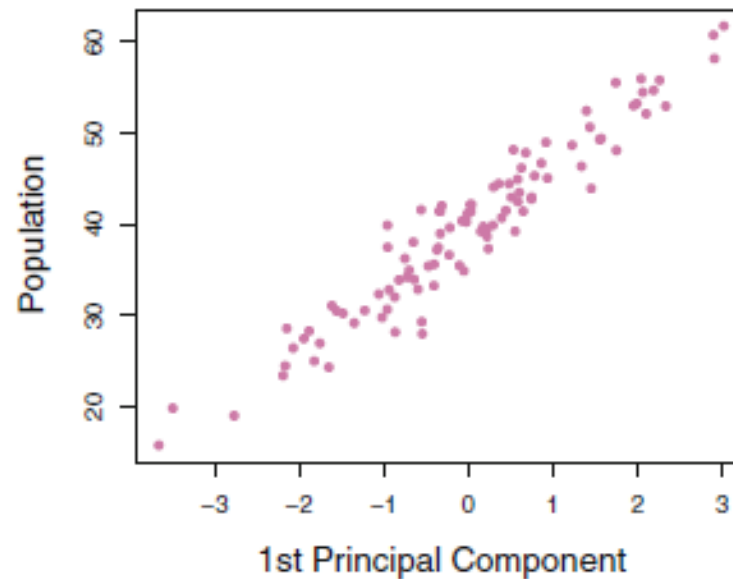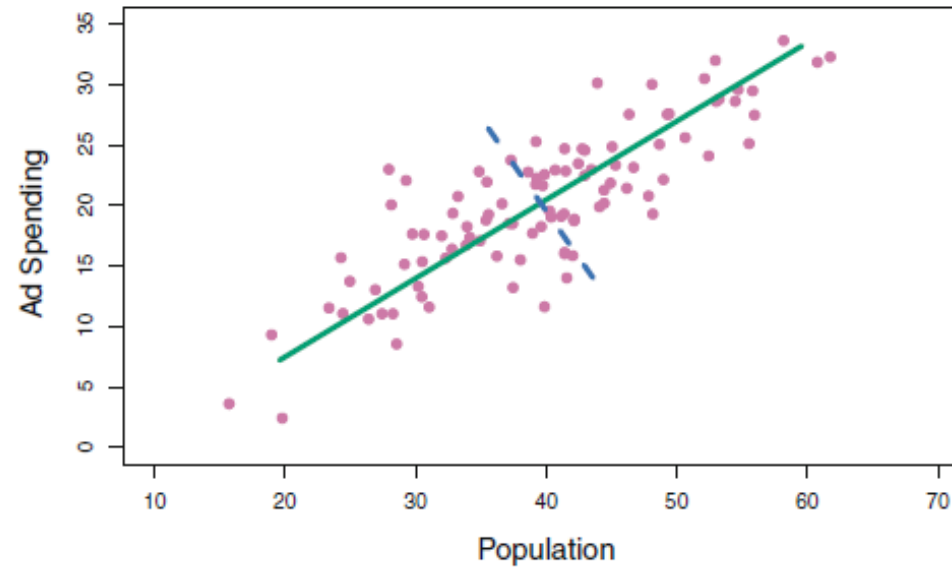
# TWO INTERPRETATIONS OF PCA

- Maximize the variance of projection along principal component

- Minimize the reconstruction error between original data and projected components
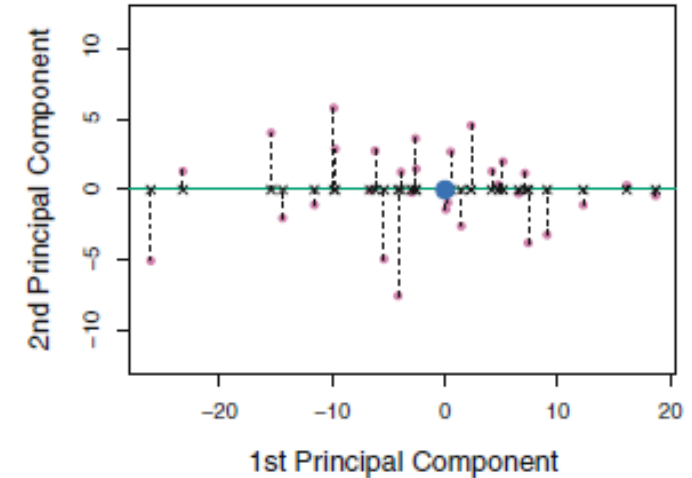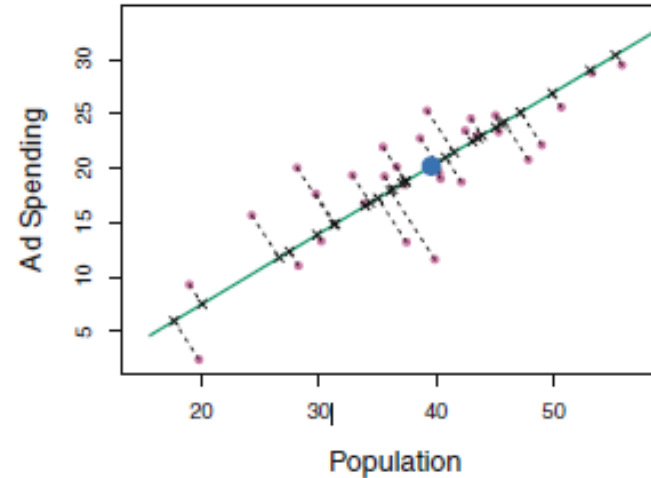


**Maximize** variance
(squared distance)
of red dots in
this direction

**Minimize** residuals
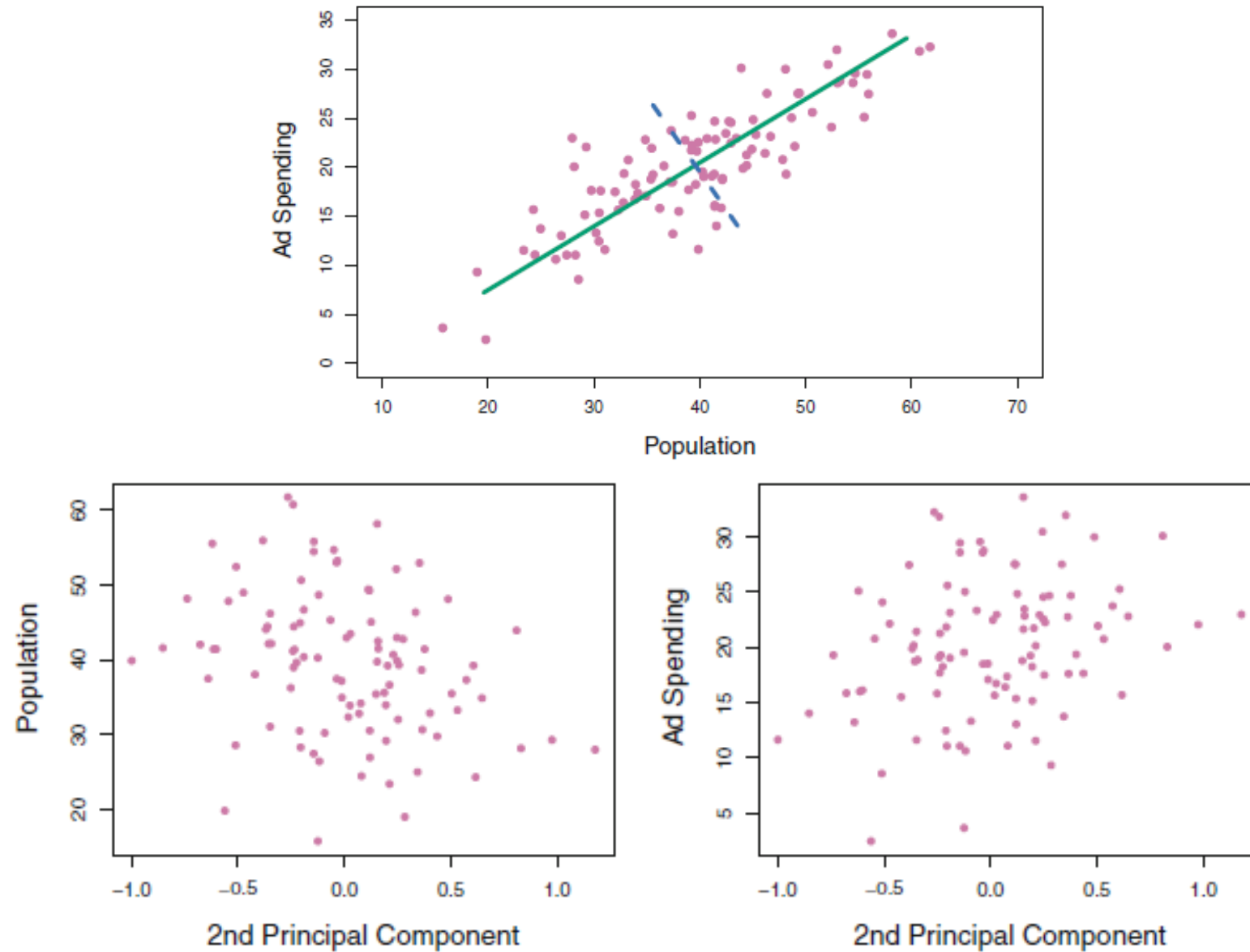(squared distance)
in this direction

# 2<sup>ND</sup> PRINCIPAL COMPONENT

- Second principal component

  - Orthogonal to the first principal component

  - And has largest variance

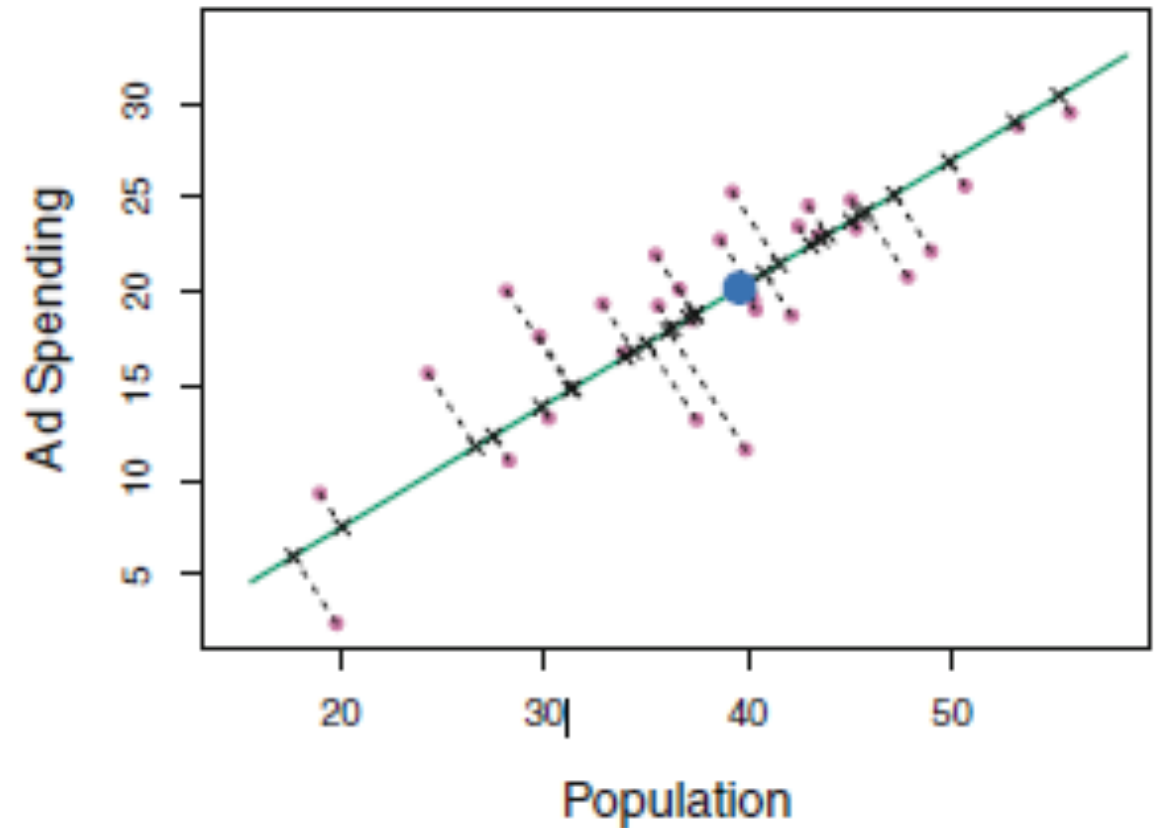- In general, we can construct up to p principal components (p features)

# PRINCIPAL COMPONENT ANALYSIS (PCA)

- Principal component loadings

  - Determines the axis

  - v = (0.839, 0.544)

- Principal component scores

  - Projected position of each point on the axis

  - Z = Xv

  (Data is centered first)



$$Z_1 = 0.839 \times (\text{pop} - \overline{\text{pop}}) + 0.544 \times (\text{ad} - \overline{\text{ad}}).$$

# PROJECTION ONTO UNIT VECTORS

- Definition of dot product:
$$\mathbf{A} \cdot \mathbf{B} = ||\mathbf{A}||_2 ||\mathbf{B}||_2 \cos\theta$$

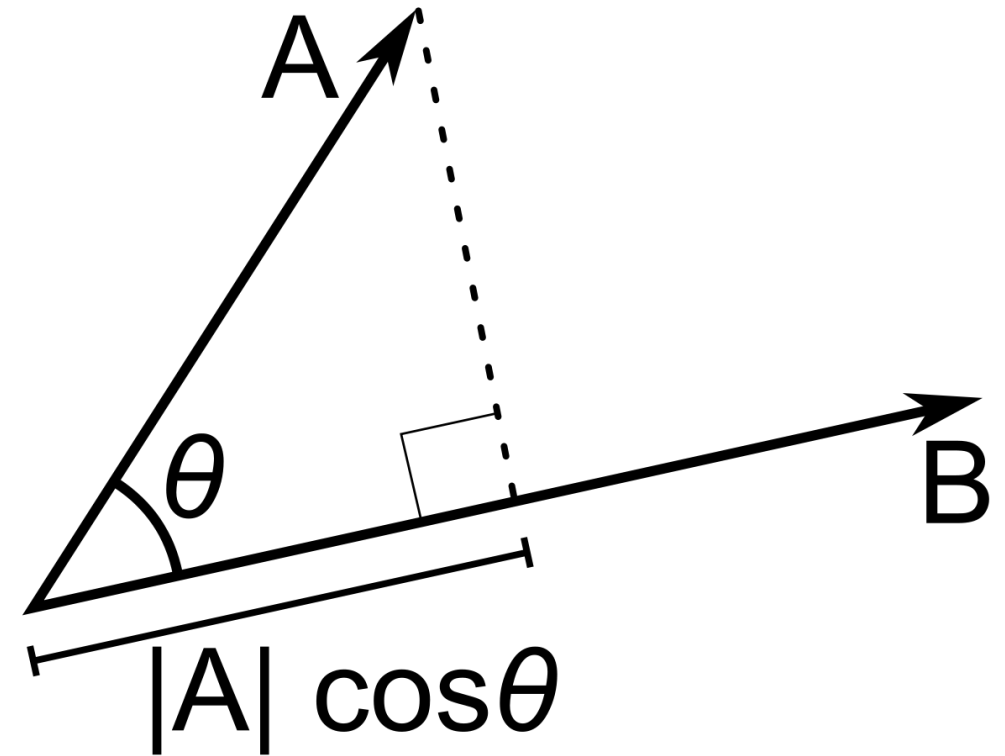- If B is a unit vector, dot product is length of the projection
$$\mathbf{A} \cdot \mathbf{B} = ||\mathbf{A}||_2 \cos\theta$$

- Projection of A onto B (vector):
$$(\mathbf{A} \cdot \mathbf{B})\mathbf{B}$$

Coefficient / score



https://en.wikipedia.org/wiki/Dot_product
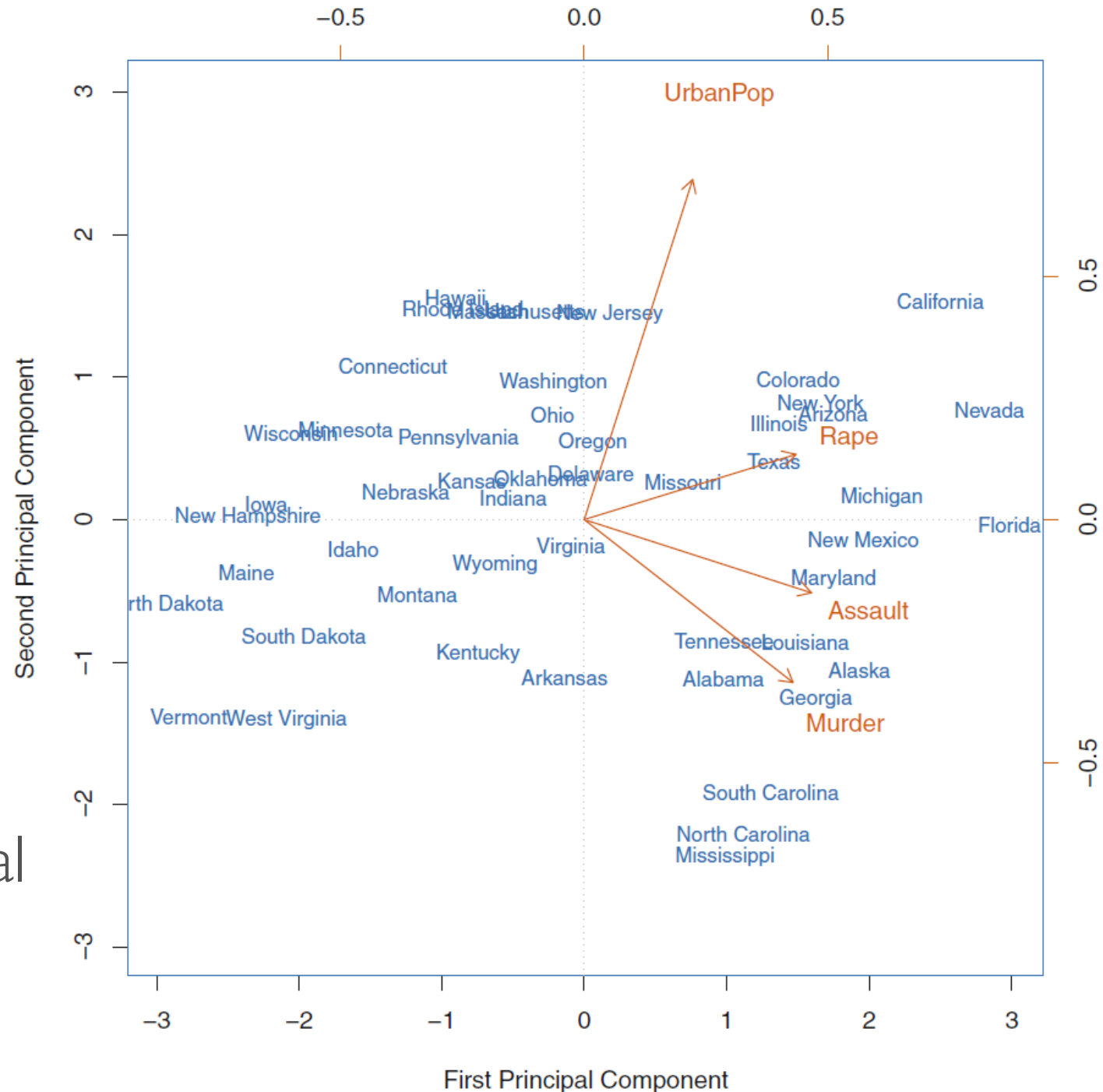
33

# EXAMPLE: US ARRESTS DATASET

- 50 states with 4 features

  - Number of arrests per 100,000 residents for Murder, Assault, Rape

  - Percent of population living in urban areas (UrbanPop)

| 1 | | Murder | Assault | UrbanPop | Rape |
|---|---|---|---|---|---|
| 2 | Alabama | 13.2 | 236 | 58 | 21.2 |
| 3 | Alaska | 10 | 263 | 48 | 44.5 |
| 4 | Arizona | 8.1 | 294 | 80 | 31 |
| 5 | Arkansas | 8.8 | 190 | 50 | 19.5 |
| 6 | California | 9 | 276 | 91 | 40.6 |
| 7 | Colorado | 7.9 | 204 | 78 | 38.7 |
| 8 | Connecticut | 3.3 | 110 | 77 | 11.1 |
| 9 | Delaware | 5.9 | 238 | 72 | 15.8 |
| 10 | Florida | 15.4 | 335 | 80 | 31.9 |
| 11 | Georgia | 17.4 | 211 | 60 | 25.8 |

# EXAMPLE: US ARRESTS (BIPLOT)

|          | PC1        | PC2         |
|----------|------------|-------------|
| Murder   | 0.5358995  | −0.4181809  |
| Assault  | 0.5831836  | −0.1879856  |
| UrbanPop | 0.2781909  | 0.8728062   |
| Rape     | 0.5434321  | 0.1673186   |

Biplot: displays both principal component scores and principal component loadings.

# PCA: 1ST PC

- 1st PC of X is unit vector that maximizes the sample variance compared to all other unit vectors

$$\mathbf{v}_1 = \operatorname{argmax}_{||\mathbf{v}||_2 = 1} (\mathbf{Xv})^\top (\mathbf{Xv})$$

- 1st PC score: $\mathbf{Xv}_1$

- Variance explained by first PC: $(\mathbf{Xv}_1)^\top (\mathbf{Xv}_1)/n$

# PCA: NEXT PC

- Idea: Successively find orthogonal directions of highest variance

- Why orthogonal?

  - Want to minimize redundancy

  - Want to look at variance in different direction

  - Computation is easier

# PCA: 2ND PC

- 2nd PC of X is unit vector that is orthogonal to the 1st PC such that it maximizes the sample variance compared to all other unit vectors that are orthogonal to the 1st PC

$$\mathbf{v}_2 = \mathrm{argmax}_{||\mathbf{v}||_2=1, \mathbf{v}^\top \mathbf{v}_1=0} (\mathbf{Xv})^\top (\mathbf{Xv})$$

- 2nd PC score: $\mathbf{Xv}_2$

- Variance explained by 2nd PC: $(\mathbf{Xv}_2)^\top (\mathbf{Xv}_2)/n$

# PCA: PROBLEM FORMULATION

- Recall 1st and 2nd PC

$$\mathbf{v}_1 = \underset{||\mathbf{v}||_2=1}{\operatorname{argmax}}(\mathbf{X}\mathbf{v})^\top(\mathbf{X}\mathbf{v})$$

$$\mathbf{v}_2 = \underset{||\mathbf{v}||_2=1, \mathbf{v}^\top\mathbf{v}_1=0}{\operatorname{argmax}}(\mathbf{X}\mathbf{v})^\top(\mathbf{X}\mathbf{v})$$

- Find orthonormal vectors that maximizes variance (assuming **X** is zero-centered)

# PCA: PROBLEM FORMULATION

- Given a feature matrix **X** with n data points, find **W** such that $\|\mathbf{W}\|_2 = 1$ and the Var(**XW**) is maximized and **W** consists of orthonormal vectors

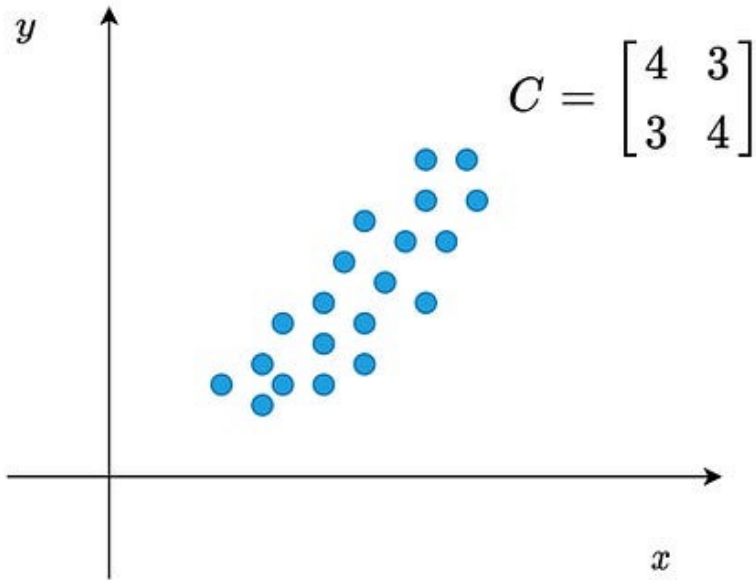$$\mathrm{Var}(\mathbf{X}\mathbf{W}) = \frac{1}{N}(\mathbf{W}^\top(\mathbf{X} - \mu_{\mathbf{X}})^\top(\mathbf{X} - \mu_{\mathbf{X}})\mathbf{W})$$

$$= \mathbf{W}^\top \boxed{\Sigma_{\mathbf{X}}} \mathbf{W}$$

Sample covariance matrix
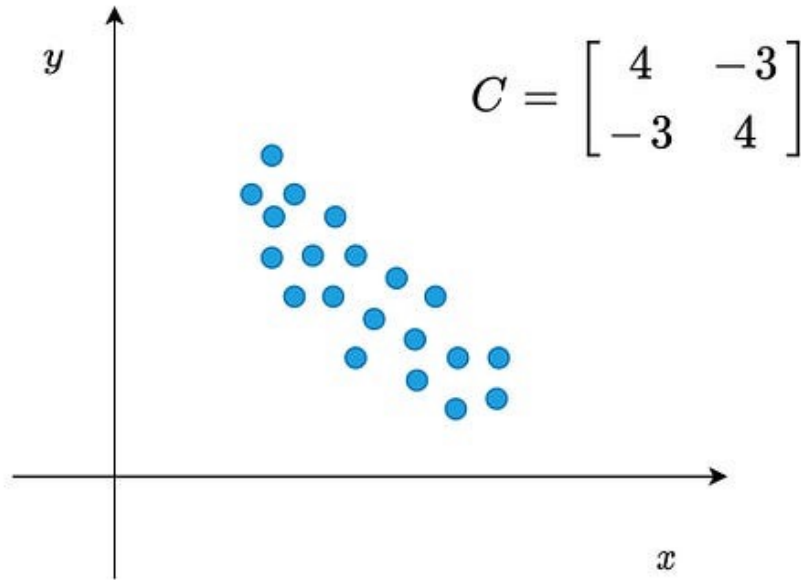
- The solution is the Eigenvectors of the covariance matrix

# EXAMPLE: COVARIANCE MATRIX

Positive
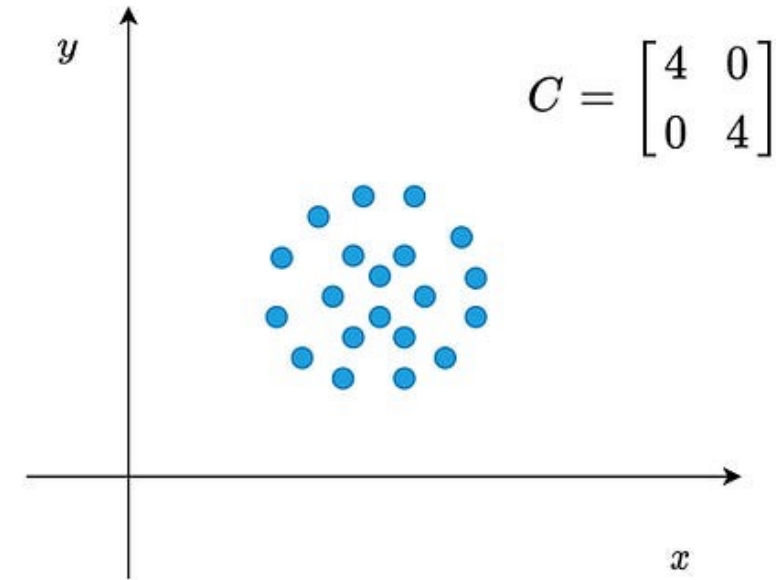Covariance

$$C = \begin{bmatrix} 4 & 3 \\ 3 & 4 \end{bmatrix}$$

Negative
Covariance

$$C = \begin{bmatrix} 4 & -3 \\ -3 & 4 \end{bmatrix}$$
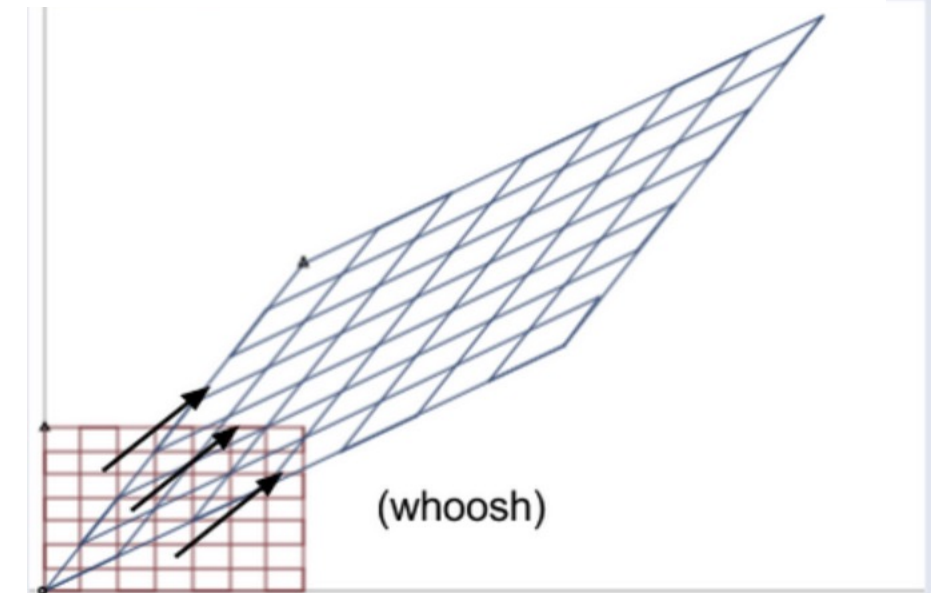
Zero
Covariance

$$C = \begin{bmatrix} 4 & 0 \\ 0 & 4 \end{bmatrix}$$
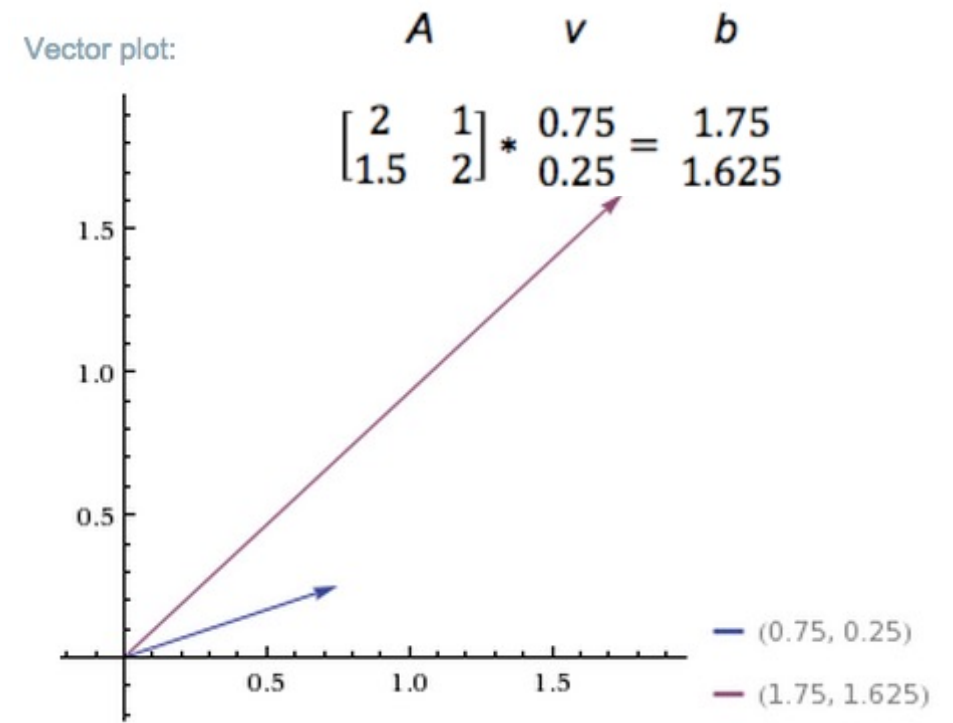
# REVIEW: EIGENVALUES + EIGENVECTORS

- Eigenvector and eigenvalue $\mathbf{Ax} = \lambda\mathbf{x}$

- Analogy: Matrix is a gust of wind (invisible force with visible result)

  - Eigenvector is like a weathervane which tells you the direction the wind is blowing in

  - Eigenvalue is just the scalar coefficient

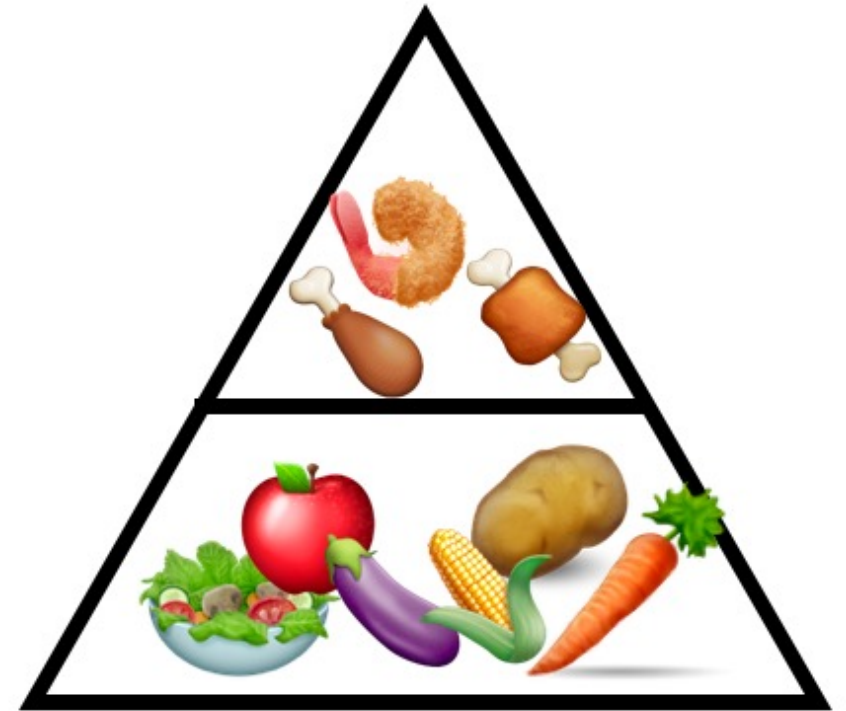- Eigenvectors of a symmetric matrix are orthonormal

Vector plot:

$$\begin{matrix} A & v & b \end{matrix}$$

$$\begin{bmatrix} 2 & 1 \\ 1.5 & 2 \end{bmatrix} * \begin{matrix} 0.75 \\ 0.25 \end{matrix} = \begin{matrix} 1.75 \\ 1.625 \end{matrix}$$

— (0.75, 0.25)
— (1.75, 1.625)

(whoosh)

# BASIC PCA ALGORITHM

- Start with a zero-centered m x n data matrix **X**

- Compute covariance matrix

- Find eigenvectors of covariance matrix

- PCs: k eigenvectors with highest eigenvalues (variance)
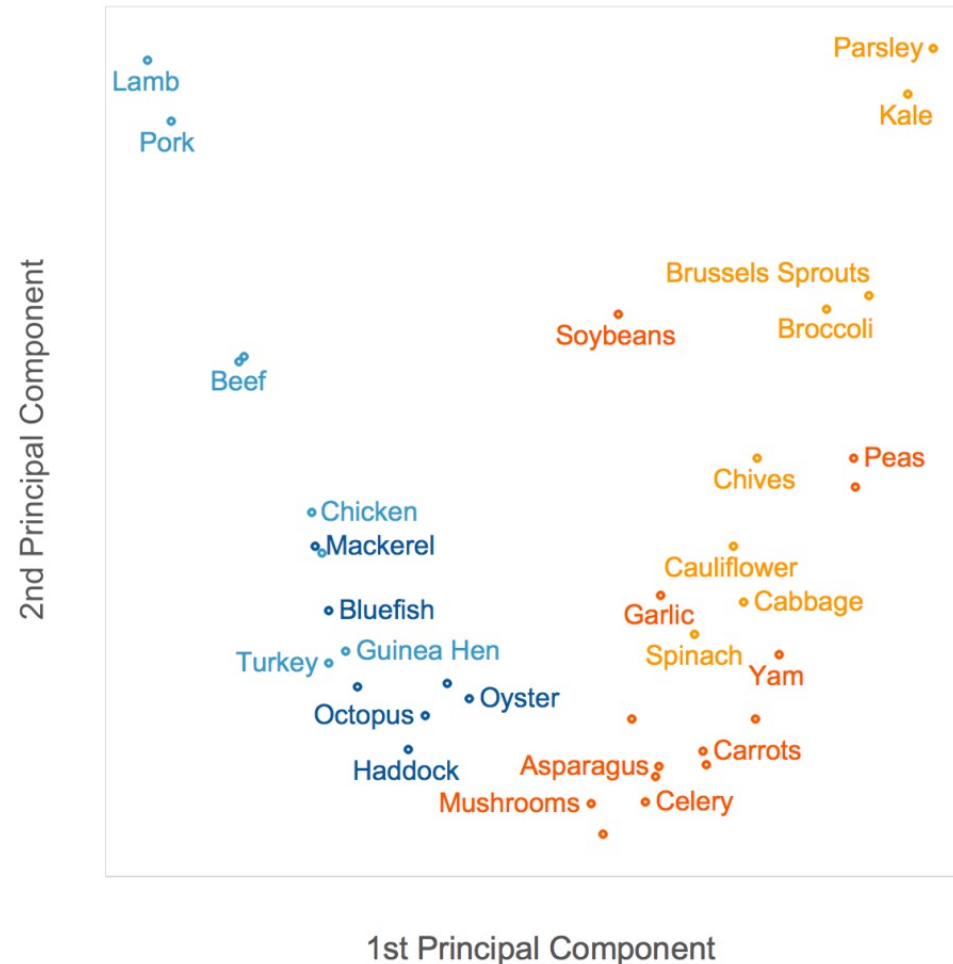
# EXAMPLE: FOOD NUTRITION

- What is the best way to differentiate food items?

  - Vitamin content

  - Protein levels

  - Fat

  - Fiber

# EXAMPLE: PCA

|           | PC1   | PC2  | PC3   | PC4   |
|-----------|-------|------|-------|-------|
| Fat       | -0.45 | 0.66 | 0.58  | 0.18  |
| Protein   | -0.55 | 0.21 | -0.46 | -0.67 |
| Fiber     | 0.55  | 0.19 | 0.43  | -0.69 |
| Vitamin C | 0.44  | 0.70 | -0.52 | 0.22  |



https://algobeans.com/2016/06/15/principal-component-analysis-tutorial/

# PCA: # OF PCS?

- How many components are sufficient to summarize the data?

# PROPORTION VARIANCE EXPLAINED

- Total variance in data (assuming zero mean):

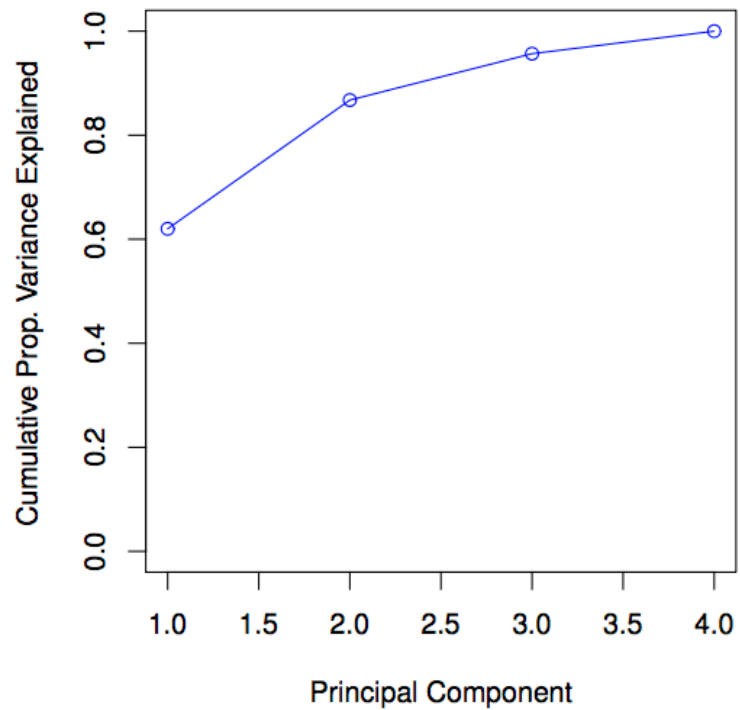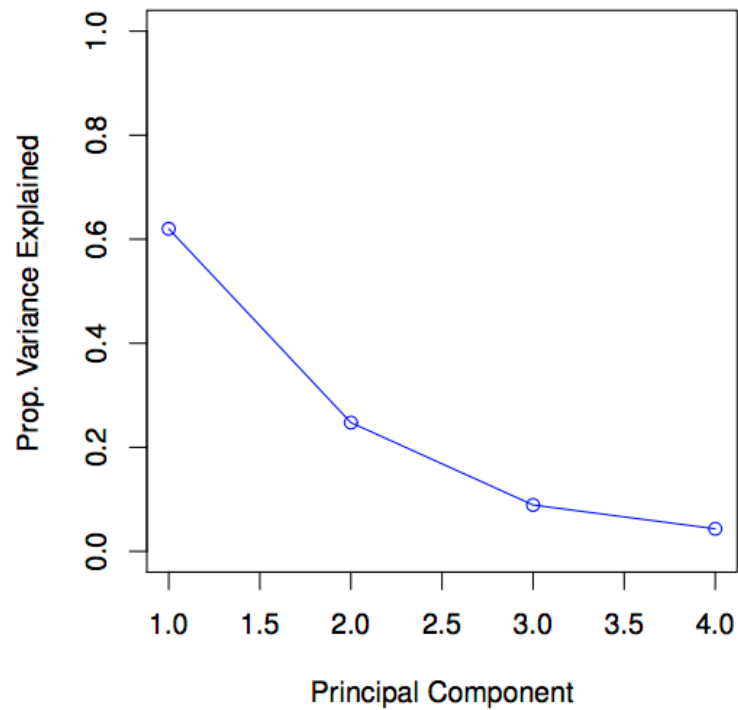$$\text{Var}(\mathbf{X}) = \sum_{j=1}^{p} \frac{1}{n} \sum_{i=1}^{n} x_{ij}^2$$

- Variance explained by the m$^{\text{th}}$ component:

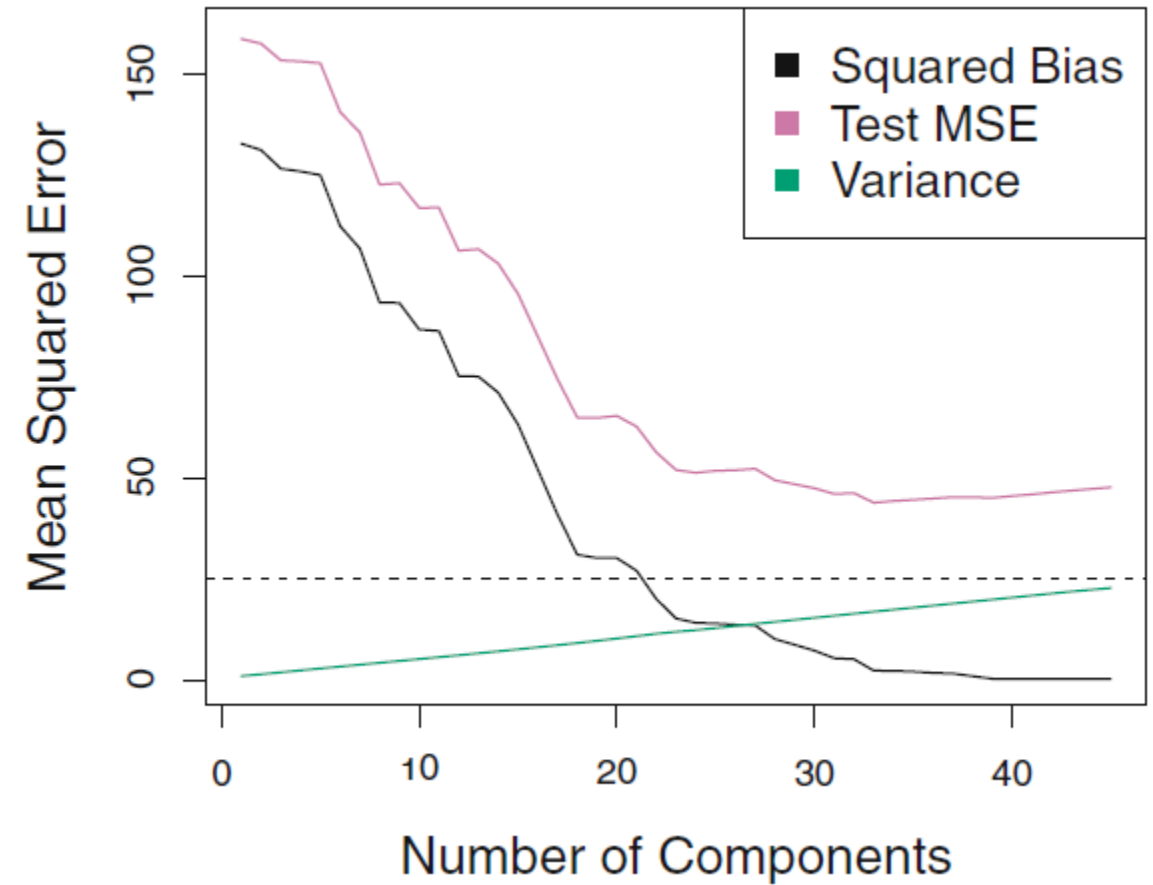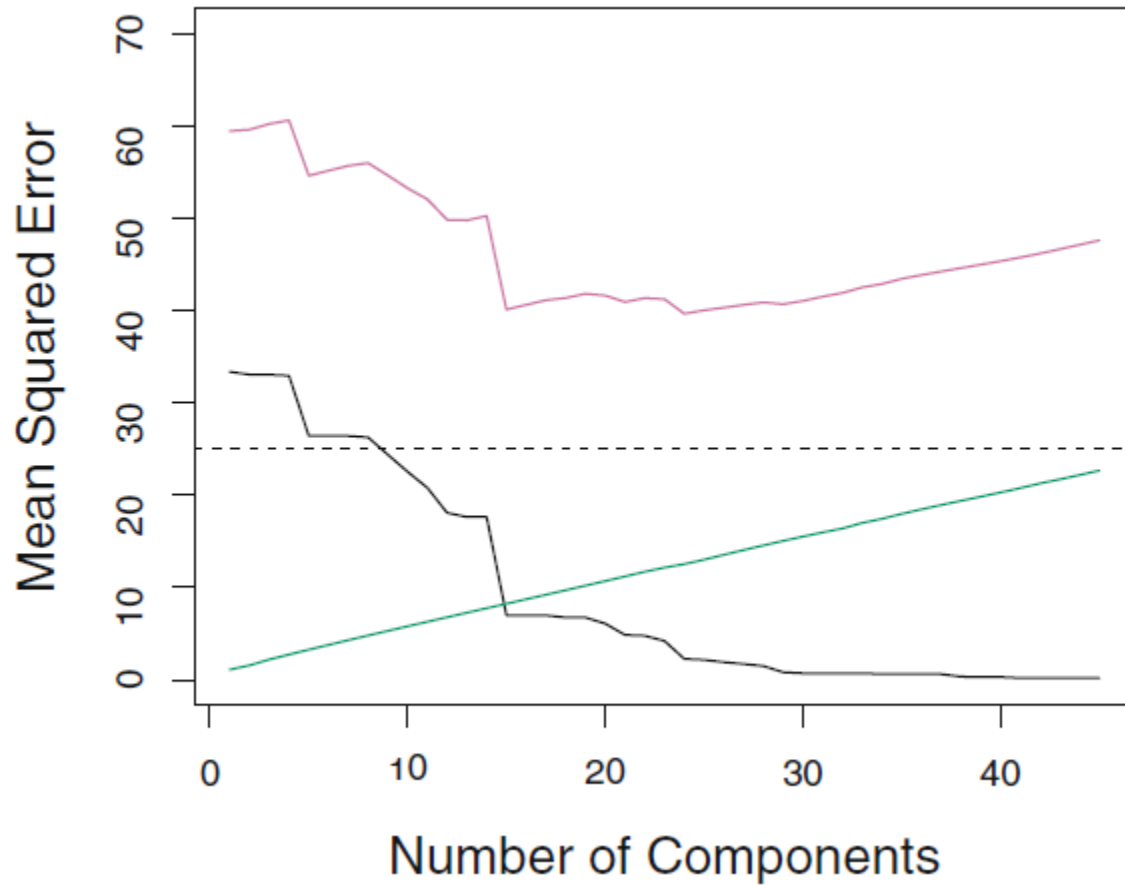$$\text{Var}(\mathbf{W}_m) = \frac{1}{n} \sum_{i=1}^{n} w_{im}^2$$

- Proportion of variance explained by m$^{\text{th}}$ component :

$$\text{PVE}_m = \frac{\text{Var}(\mathbf{W}_m)}{\text{Var}(\mathbf{X})}$$
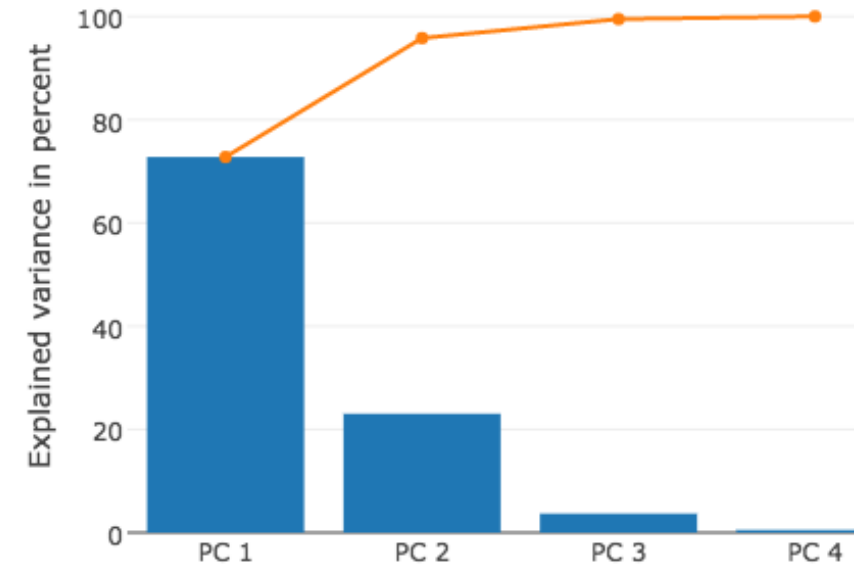
# PCA: SCREE PLOT (UNSUPERVISED)

# USE PCA FOR SUPERVISED LEARNING

# PCA: INTERPRETATION

- If variances of PCs drop off quickly, then X is highly collinear

- Reduce dimensionality of data by keeping only the PCs with highest variance



https://plot.ly/ipython-notebooks/principal-component-analysis/

# PCA: SKLEARN

```python
from sklearn.decomposition import PCA as sklearnPCA
sklearn_pca = sklearnPCA(n_components=2)
Y_sklearn = sklearn_pca.fit_transform(X_std)
```

# HOMEWORK #5

- Due **11/16 @ 11:59 PM ET** on Gradescope

- 2 questions

  - PCA

  - Almost Random Forest

# DEMO: PCA-EXAMPLE.IPYNB

HTTPS://COLAB.RESEARCH.GOOGLE.COM/DRIVE/12VDO-JD0PVFLVZZLI2FT50DYESJ8V6WW