

10-605/805 – ML for Large Datasets

Lecture 3: Principal Component Analysis

Henry Chai

9/6/22

Front Matter

- HW1 released 8/30, due 9/14 at 11:59 PM
 - **For HW1 only, the programming part is optional (but strongly encouraged)**
 - The written part is nominally about PCA but can be solved using pre-requisite knowledge (linear algebra)
- Recitations on Friday, 11:50 – 1:10 (**different from lecture**) in GHC 4401 (**same as lecture**)
 - Recitation 2 on 9/9: Review of linear algebra

Data Pre-processing

- ETL (extract-transfer-load)
- Cleaning data
 - Missing features/labels
 - Duplicated observations
 - Formatting errors
- Understanding data
 - Summarization
 - Exploration
 - Visualization

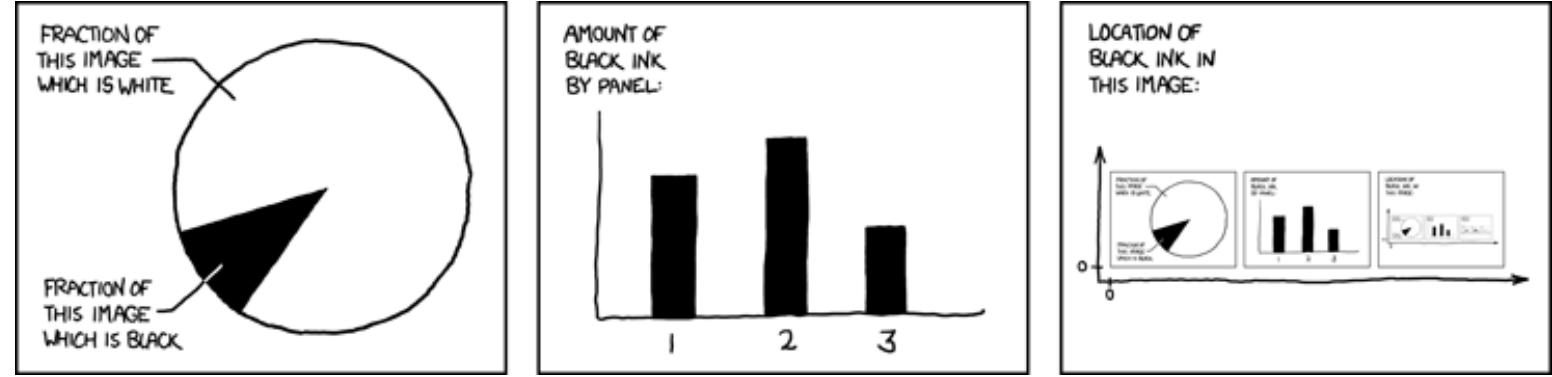
Data Pre-processing

- ETL (extract-transfer-load)
- Cleaning data
 - Missing features/labels
 - Duplicated observations
 - Formatting errors
- Understanding data
 - Summarization
 - Exploration
 - Visualization

Given some
(labelled)
dataset, what
questions can
you ask to
better
understand
the data?

- what are the units?
- what are the features?
 - and what type are they all?
- how many observations are we working w/?
- Is it complete and collected uniformly/
systematically?
- are features related to one another?
- what are actually doing?

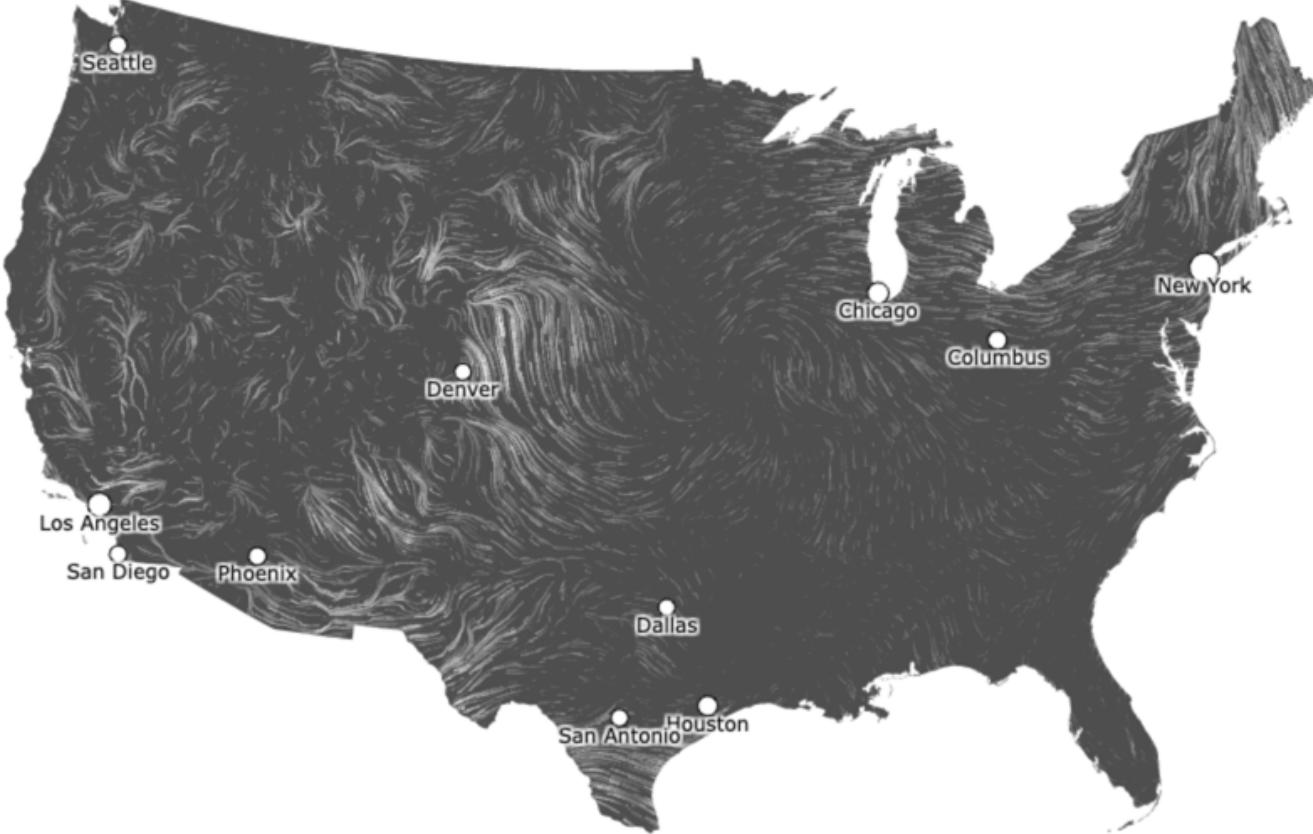
Data Visualization



- Visualizations can be used to
 - Provide insight about trends/groups/relationships
 - Reveal systematic errors
 - Aid in model selection
 - Evaluate training (e.g., measure convergence)
 - Interpret/explain predictions

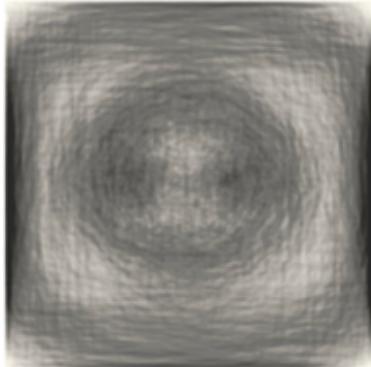
Data Visualization: Examples

- Understanding scale



Data Visualization: Examples

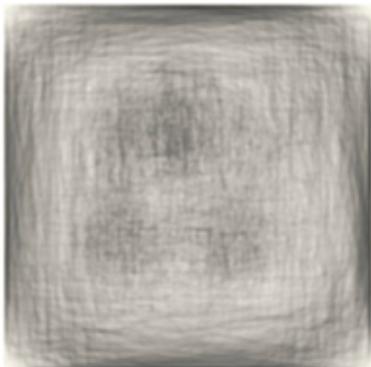
- Understanding clusters



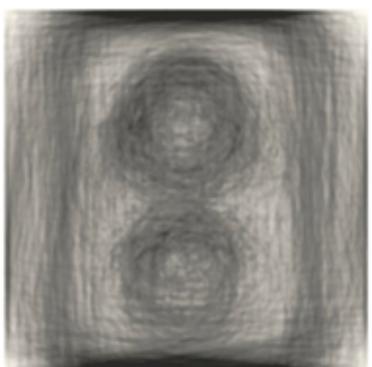
Germany



Japan



Malaysia

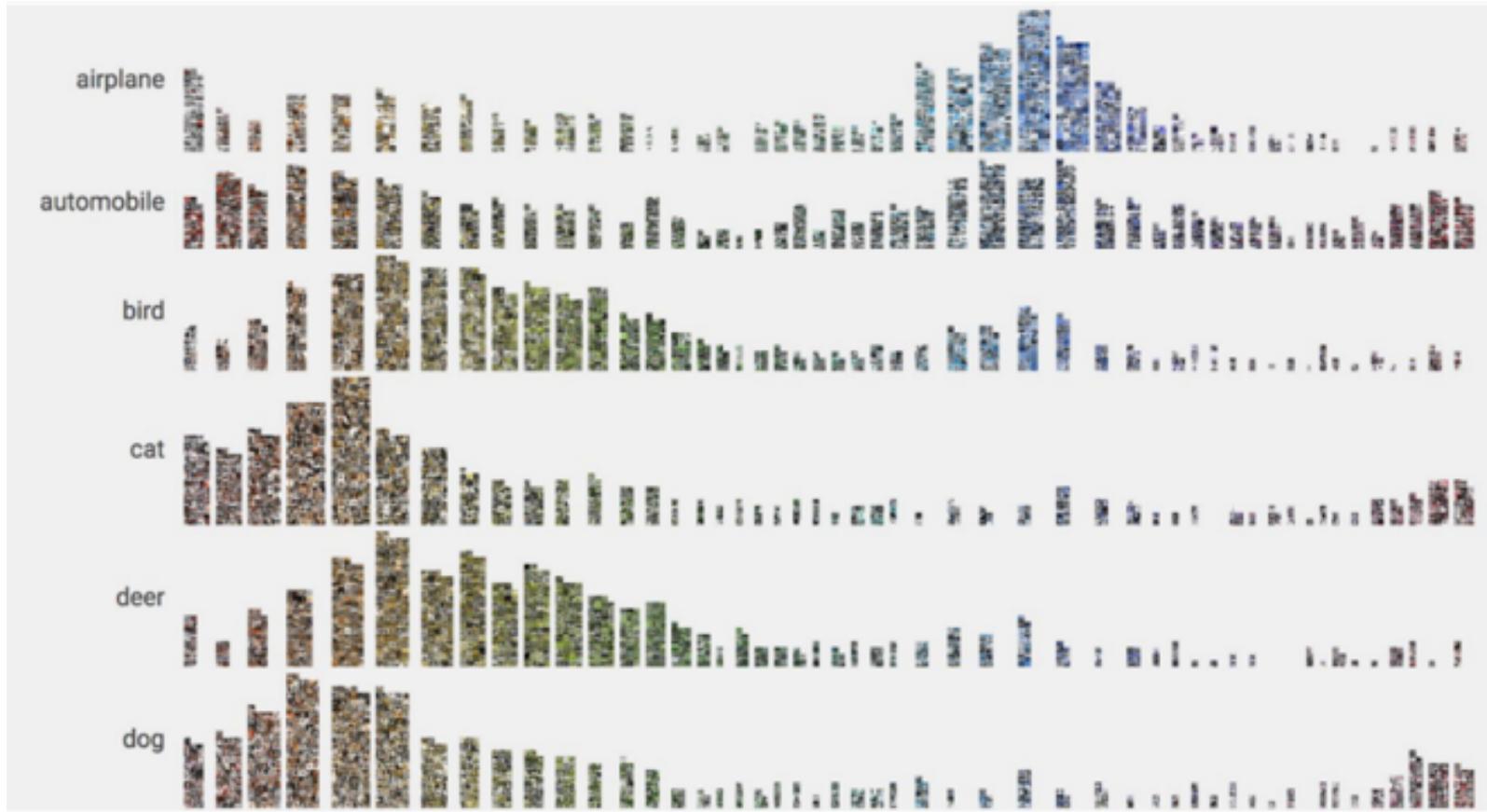


Sweden

electrical
outlets

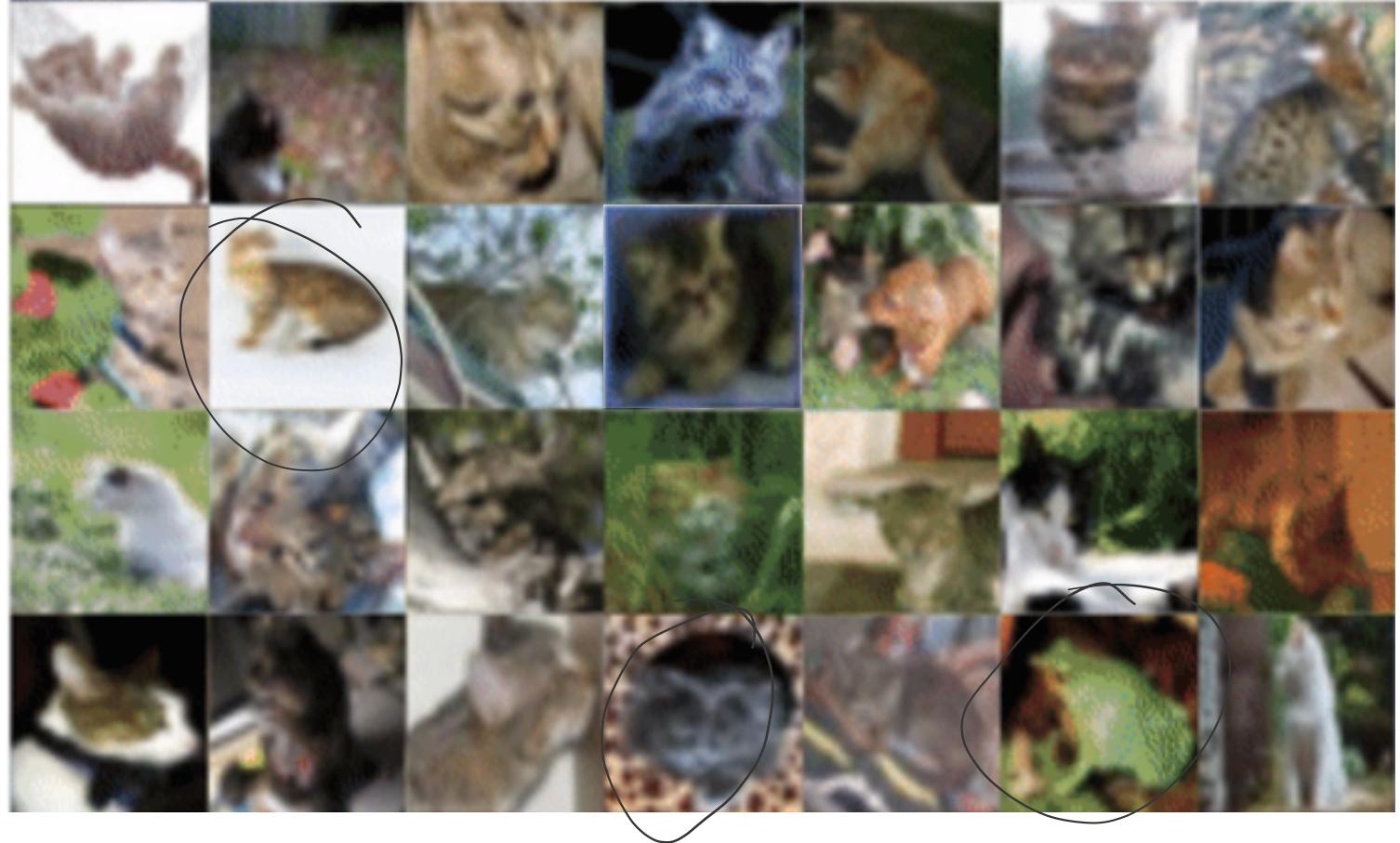
Data Visualization: Examples

- Understanding variation



Data Visualization: Examples

- Understanding mistakes



Common Data Visualizations

- Summary statistics
- Box plots
- Histograms
- Scatter plots

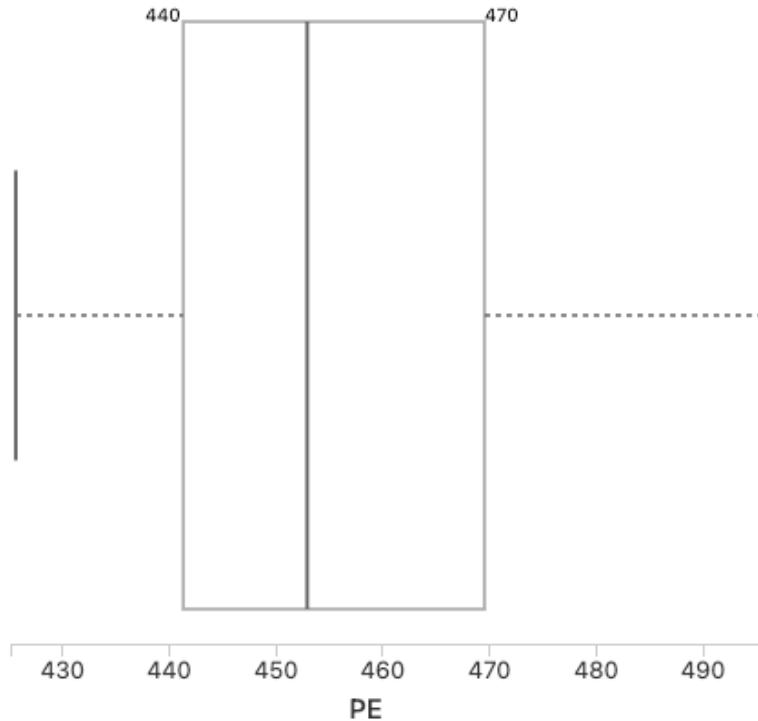
Common Data Visualizations

- Summary statistics
- Box plots
- Histograms
- Scatter plots

	AT	PE
Mean	19.65	454.37
Stddev	7.45	17.07
Min	1.81	420.26
Max	37.11	495.76

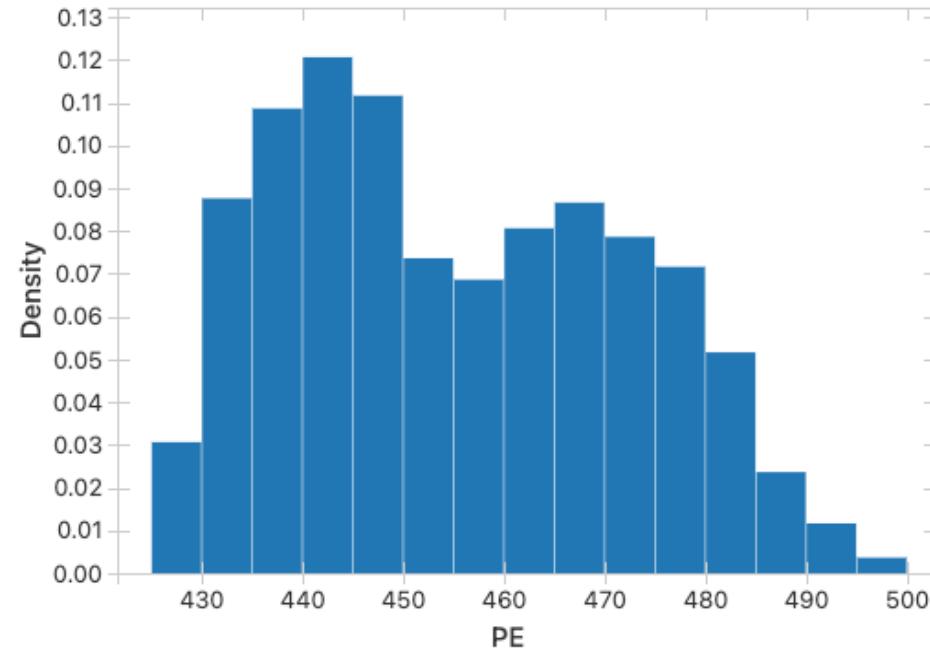
Common Data Visualizations

- Summary statistics
- Box plots
- Histograms
- Scatter plots



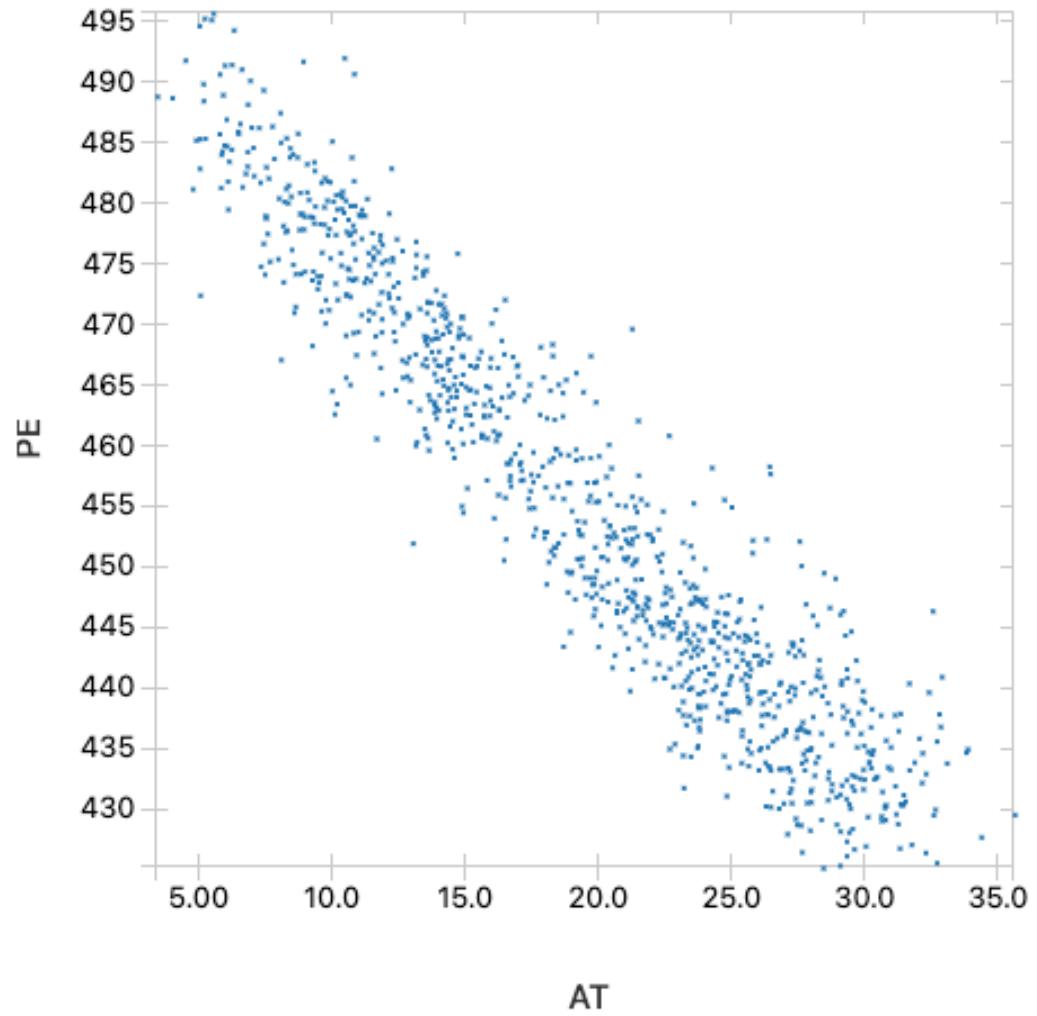
Common Data Visualizations

- Summary statistics
- Box plots
- Histograms
- Scatter plots



Common Data Visualizations

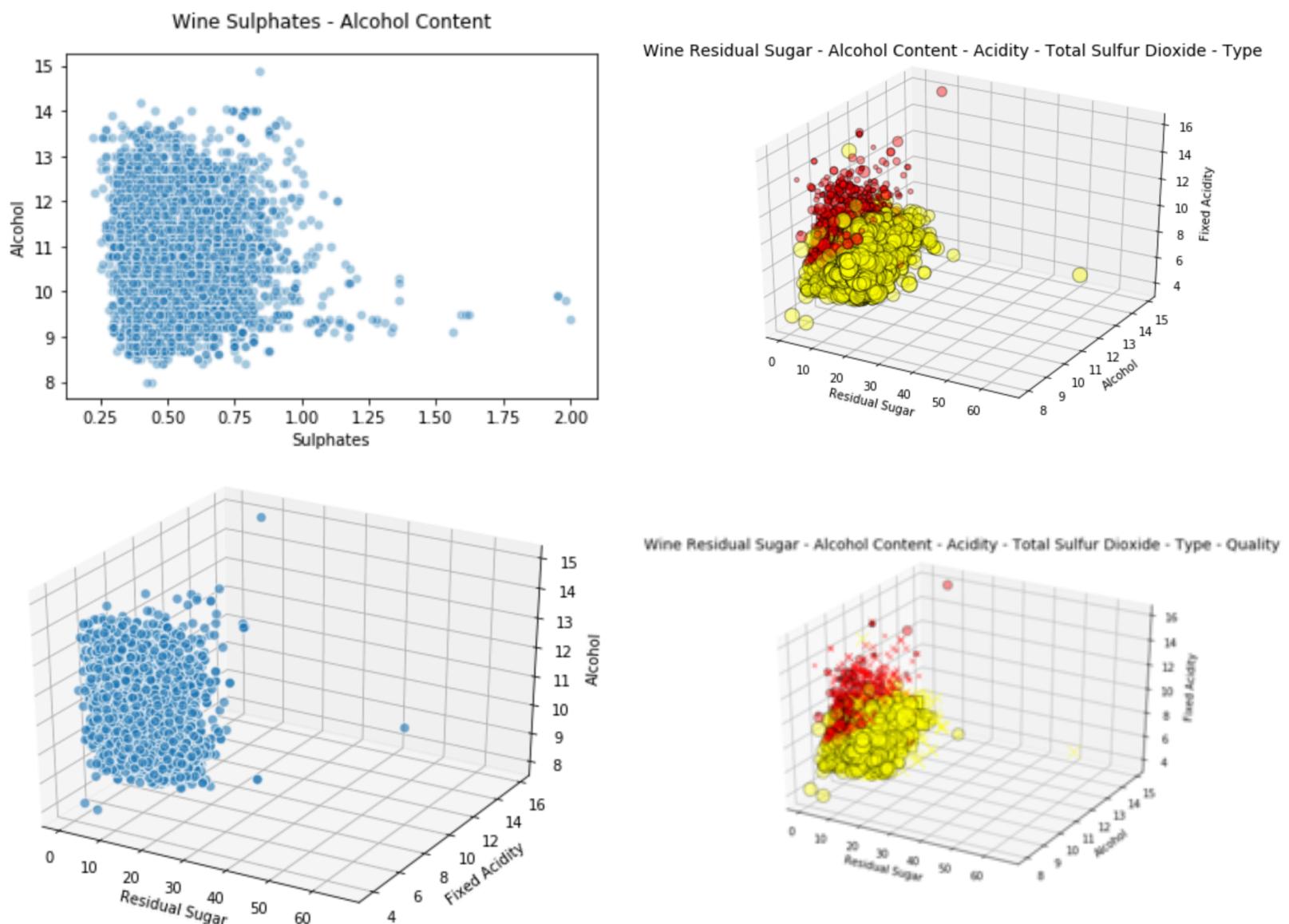
- Summary statistics
- Box plots
- Histograms
- Scatter plots



Big Data Visualizations

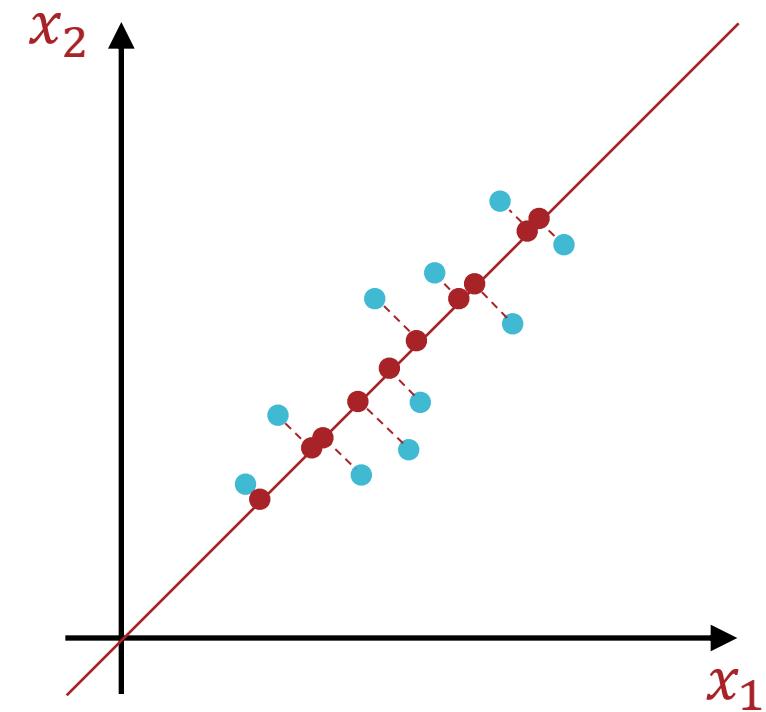
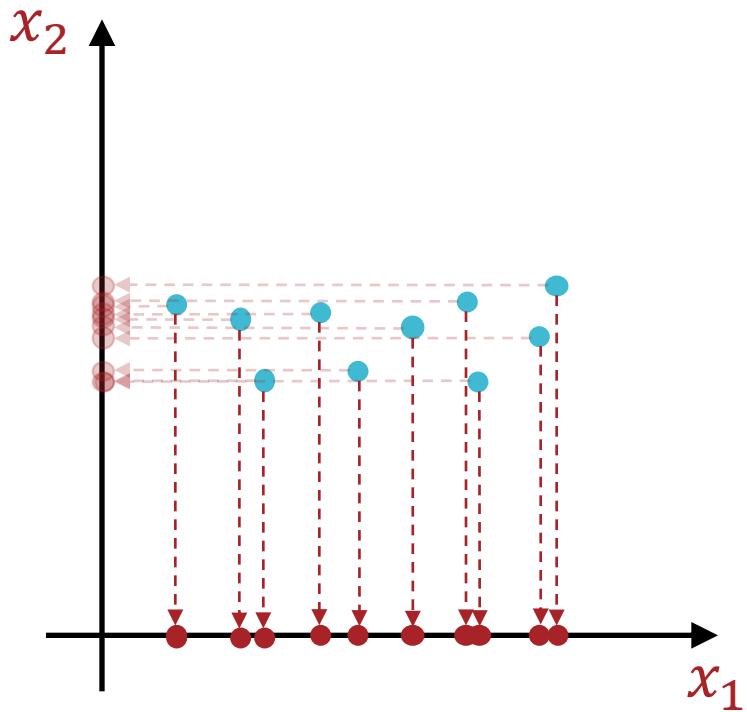
- Large n
 - Computationally expensive to render
 - Dense/complex
 - Address via subsampling or parallelization
- Large k
 - Difficult to represent more than a few dimensions

Big Data Visualizations

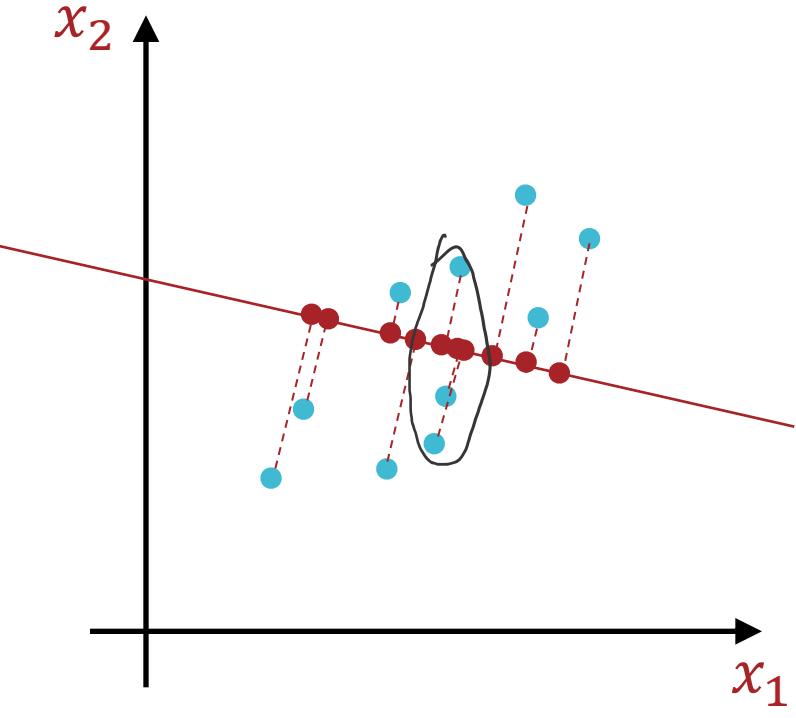


Big Data Visualizations

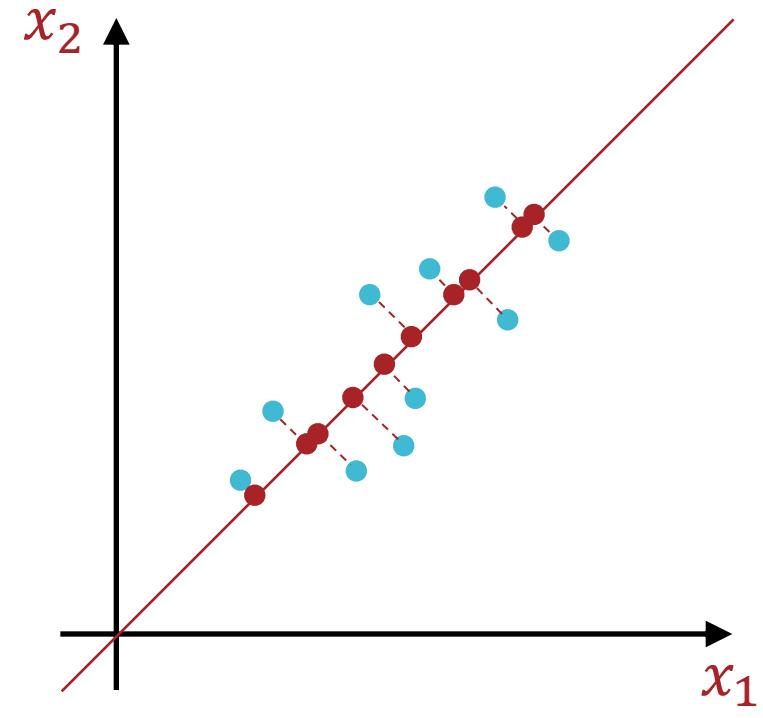
- Large n
 - Computationally expensive to render
 - Dense/complex
 - Address via subsampling or parallelization
- Large k
 - Difficult to represent more than a few dimensions
 - Address via dimensionality reduction = learning a latent (typically lower-dimensional) representation



Feature Elimination

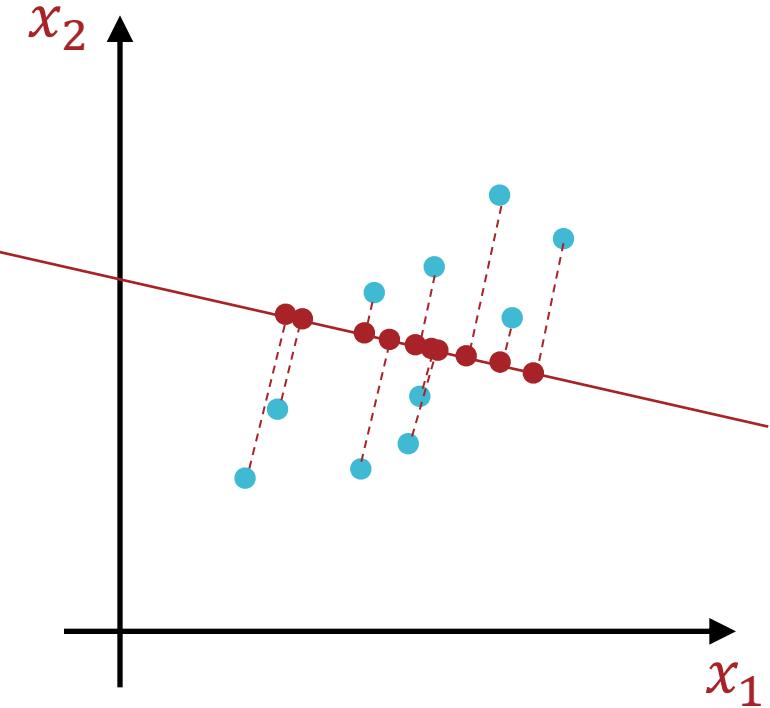


Option A

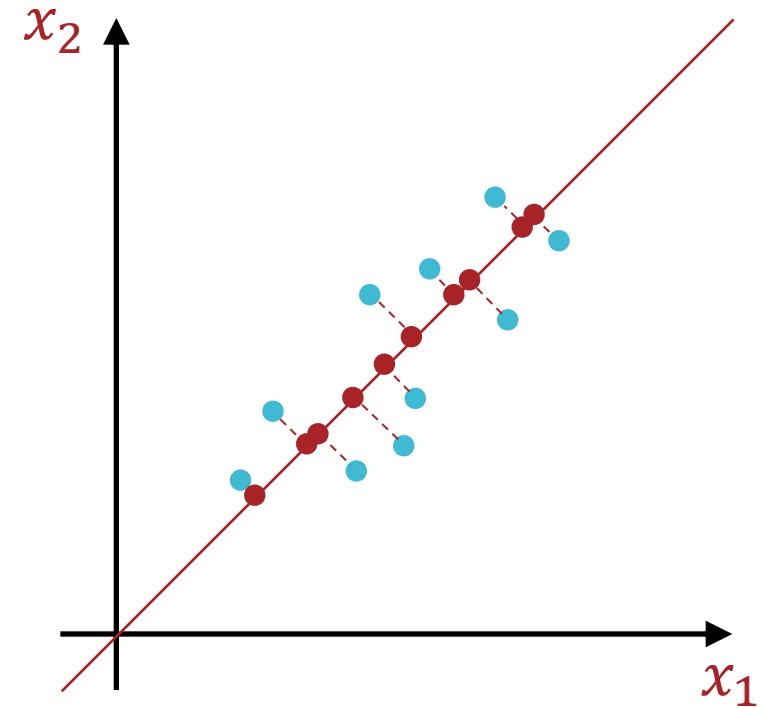


Option B

Dimensionality Reduction

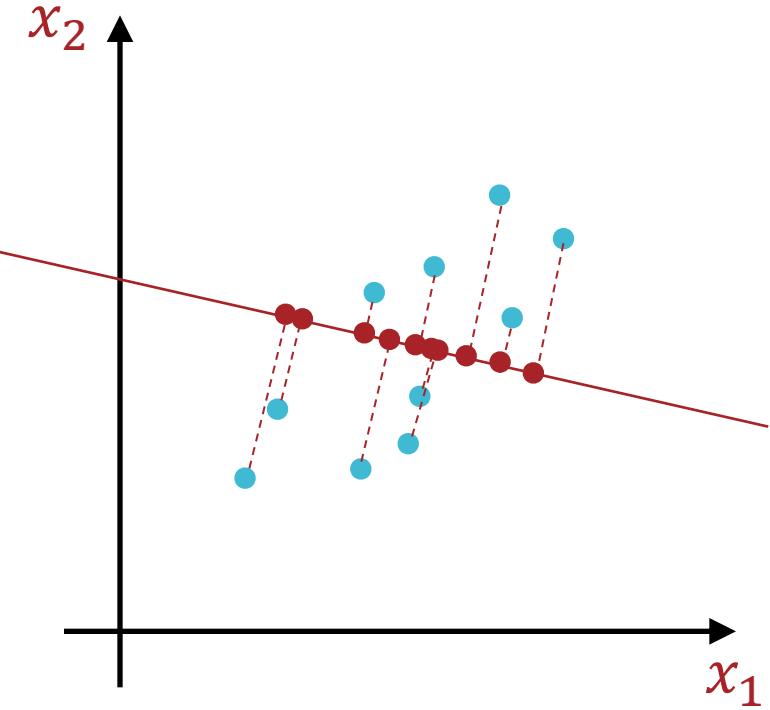


Option A

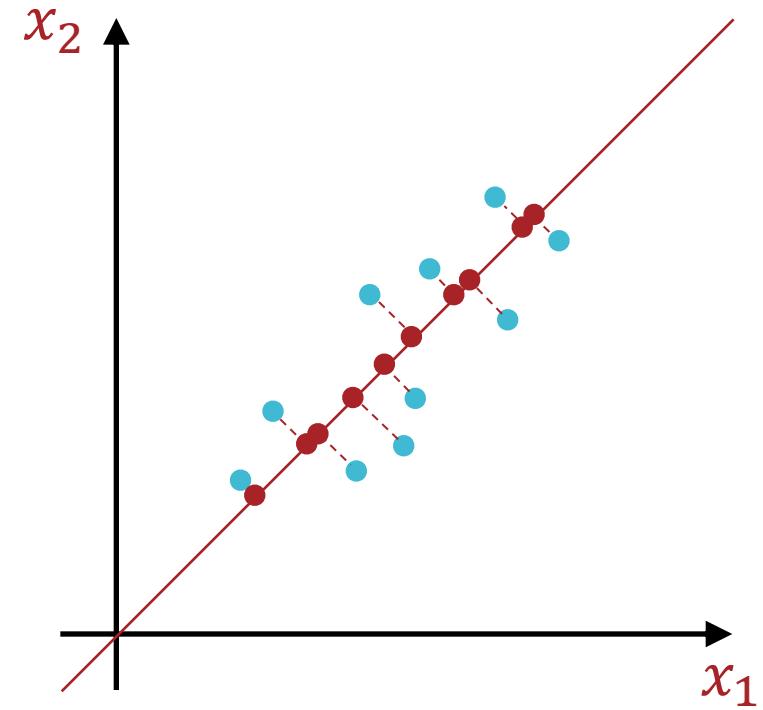


Option B

Which projection do you prefer?

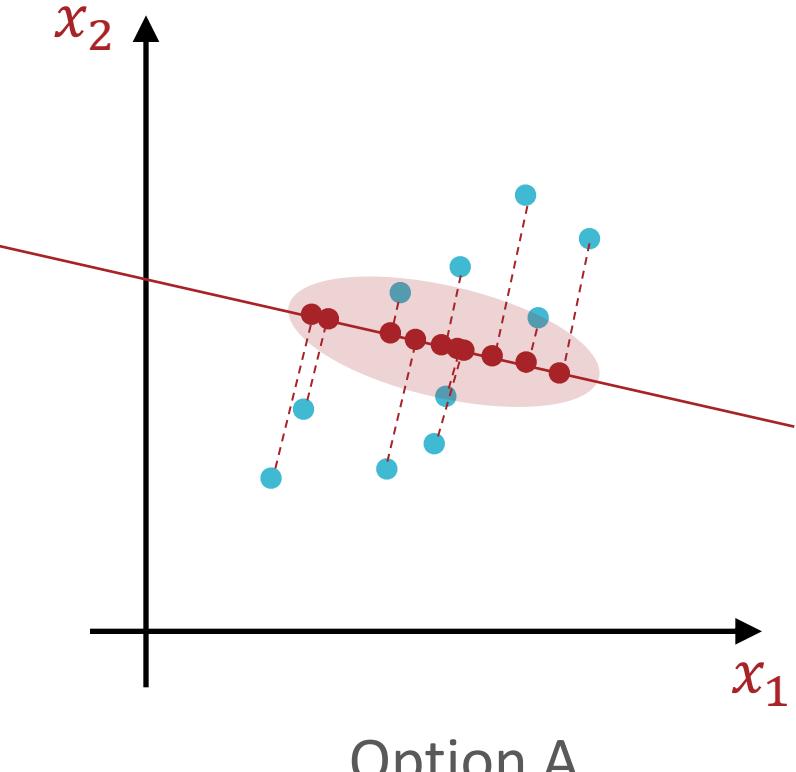


Option A

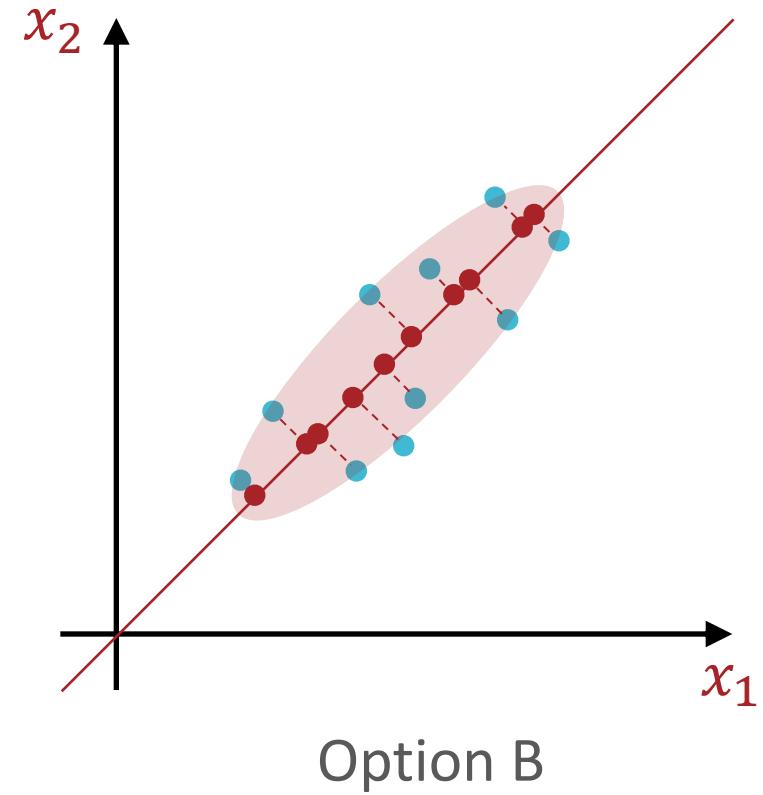


Option B

Goal: minimize the reconstruction error



Option A



Option B

Goal: maximize the variance of the projections

Centering the Data

- To be consistent, we will constrain principal components to be *orthonormal vectors* (orthogonal unit vectors) that begin at the origin
- Preprocess data to be centered around the origin:

$$1. \boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \boldsymbol{x}^{(i)}$$

$$2. \tilde{\boldsymbol{x}}^{(i)} = \boldsymbol{x}^{(i)} - \boldsymbol{\mu} \quad \forall i$$

$$3. X = \underbrace{\begin{bmatrix} \tilde{\boldsymbol{x}}^{(1)^T} \\ \tilde{\boldsymbol{x}}^{(2)^T} \\ \vdots \\ \tilde{\boldsymbol{x}}^{(n)^T} \end{bmatrix}}_{\in \mathbb{R}^{n \times k}}$$

Reconstruction Error

- The projection of $\tilde{\mathbf{x}}^{(i)}$ onto a vector \mathbf{v} is

$$\mathbf{z}^{(i)} = \left(\frac{\mathbf{v}^T \tilde{\mathbf{x}}^{(i)}}{\|\mathbf{v}\|_2} \right) \frac{\mathbf{v}}{\|\mathbf{v}\|_2}$$

Length of projection

Direction of projection

Reconstruction Error

- The projection of $\tilde{x}^{(i)}$ onto a unit vector v is

$$z^{(i)} = (v^T \tilde{x}^{(i)}) v$$

$$\begin{aligned}
 \hat{v} &= \underset{\|v\|_2=1}{\operatorname{argmin}} \quad (\text{min reconstruction error}) \\
 &= \underset{\|v\|_2=1}{\operatorname{argmin}} \sum_{i=1}^n \|\tilde{x}^{(i)} - z^{(i)}\|_2^2 \\
 &= \underset{\|v\|_2=1}{\operatorname{argmin}} \sum_{i=1}^n \|\tilde{x}^{(i)} - (v^T \tilde{x}^{(i)}) v\|_2^2 \\
 \|\tilde{x}^{(i)} - (v^T \tilde{x}^{(i)}) v\|_2^2 &= (\tilde{x}^{(i)T} \tilde{x}^{(i)} - 2v^T \tilde{x}^{(i)} (v^T \tilde{x}^{(i)}) \\
 &\quad + v^T \tilde{x}^{(i)} (v^T \tilde{x}^{(i)}) v) \underbrace{\|v\|_2^2}_{} = 1
 \end{aligned}$$

Minimizing the Reconstruction Error

$$\hat{v} = \operatorname{argmin}_{v: \|v\|_2^2=1} \sum_{i=1}^n \|\tilde{x}^{(i)}\|_2^2 - (v^T \tilde{x}^{(i)})^2$$

$$E[\sqrt{v^T \tilde{x}^{(i)}}] = \sigma$$

$$\begin{aligned} \hat{v} &= \operatorname{argmin}_{v: \|v\|_2^2=1} \sum_{i=1}^n -(\sqrt{v^T \tilde{x}^{(i)}})^2 = \operatorname{argmax}_{v: \|v\|_2^2=1} \sum_{i=1}^n (\sqrt{v^T \tilde{x}^{(i)}})^2 \\ &= \operatorname{argmax}_{v: \|v\|_2^2=1} \sum_{i=1}^n (v^T \tilde{x}^{(i)} \tilde{x}^{(i)T} v) = \operatorname{argmax}_{v: \|v\|_2^2=1} v^T \left(\sum_{i=1}^n \tilde{x}^{(i)} \tilde{x}^{(i)T} \right) v \end{aligned}$$

$$= \operatorname{argmax}_{v: \|v\|_2^2=1} v^T X^T X v$$

$$= \operatorname{argmax}_{v: \|v\|_2^2=1} v^T C_X v$$

Covariance matrix

Maximizing the Variance

$$\hat{v} = \underset{\|v\|_2^2=1}{\operatorname{argmax}} v^T C_X v$$

$$\begin{aligned} L(v, \lambda) &= v^T C_X v - \lambda(\|v\|_2^2 - 1) \\ &= \underbrace{v^T C_X v}_{\frac{\partial L}{\partial v}} - \lambda(\underbrace{v^T v - 1}_{\frac{\partial L}{\partial \lambda}}) \end{aligned}$$

$$\begin{aligned} \frac{\partial L}{\partial v} &= C_X v - \lambda v \\ \Rightarrow C_X \hat{v} - \lambda \hat{v} &= 0 \Rightarrow C_X \hat{v} = \lambda \hat{v} \end{aligned}$$

λ is an eigenvalue of C_X and
its corresponding eigenvector \hat{v}

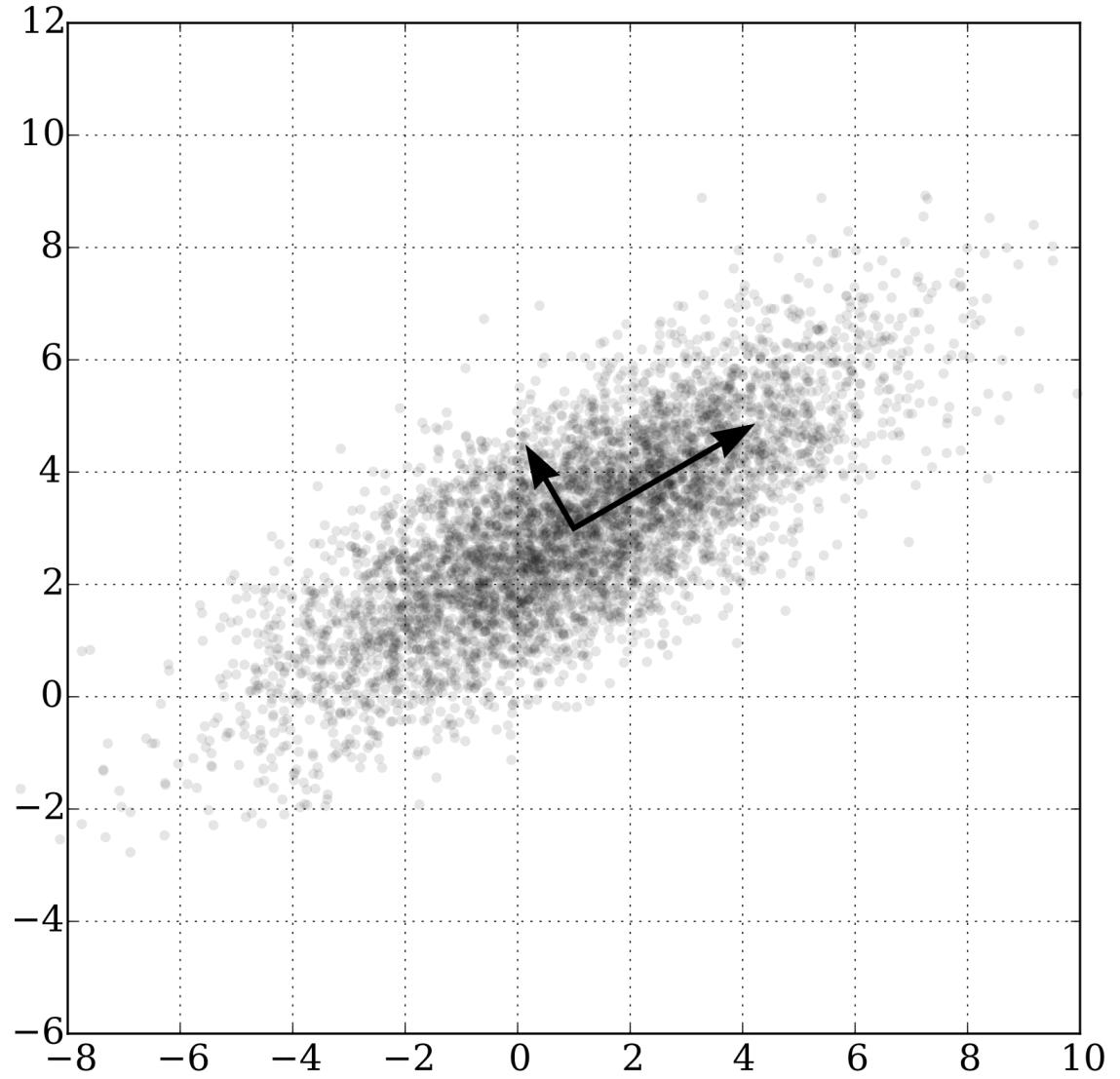
Maximizing the Variance

$$\hat{\mathbf{v}} = \underset{\mathbf{v}: \|\mathbf{v}\|_2^2=1}{\operatorname{argmax}} \mathbf{v}^T C_X \mathbf{v}$$

$$C_X \hat{\mathbf{v}} = \lambda \hat{\mathbf{v}} \rightarrow \hat{\mathbf{v}}^T C_X \hat{\mathbf{v}} = \lambda \hat{\mathbf{v}}^T \hat{\mathbf{v}} = \lambda$$

- The first principal component is the eigenvector $\hat{\mathbf{v}}_1$ that corresponds to the largest eigenvalue λ_1
- The second principal component is the eigenvector $\hat{\mathbf{v}}_2$ that corresponds to the second largest eigenvalue λ_2
 - $\hat{\mathbf{v}}_1$ and $\hat{\mathbf{v}}_2$ are orthonormal!
- Etc ...
- λ_i is a measure of how much variance falls along $\hat{\mathbf{v}}_i$

Principal Components: Example



HW1 Aside: Sparse PCA

$$\hat{\boldsymbol{v}} = \operatorname{argmax}_{\boldsymbol{v}: \|\boldsymbol{v}\|_2^2=1 \text{ and } \|\boldsymbol{v}\|_0 \leq s} \boldsymbol{v}^T C_X \boldsymbol{v}$$

- L0-norm constraint: \boldsymbol{v} can contain at most s non-zero elements
- Reduces to standard PCA if $s = k$
- Sparse principal components may be easier to interpret and can also reduce data pre-processing needs

PCA Algorithm

- Input: $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^n, r$
 1. Center the data
 - A. Optionally, normalize the data by features so that all features are of the same scale
 2. Compute the covariance matrix $C_X = X^T X$
 3. Collect the top r eigenvectors (corresponding to the r largest eigenvalues), $P \in \mathbb{R}^{k \times r}$
 4. Project the data into the space defined by P , $Z = XP$
- Output: Z , the latent representation (“PCA scores”)

How many PCs should we use?

- Input: $\mathcal{D} = \{(\mathbf{x}^{(i)})\}_{i=1}^n$, r
- 1. Center the data
 - A. Optionally, normalize the data by features so that all features are of the same scale
- 2. Compute the covariance matrix $C_X = X^T X$
- 3. Collect the top r eigenvectors (corresponding to the r largest eigenvalues), $P \in \mathbb{R}^{k \times r}$
- 4. Project the data into the space defined by P , $Z = XP$
- Output: Z , the latent representation (“PCA scores”)

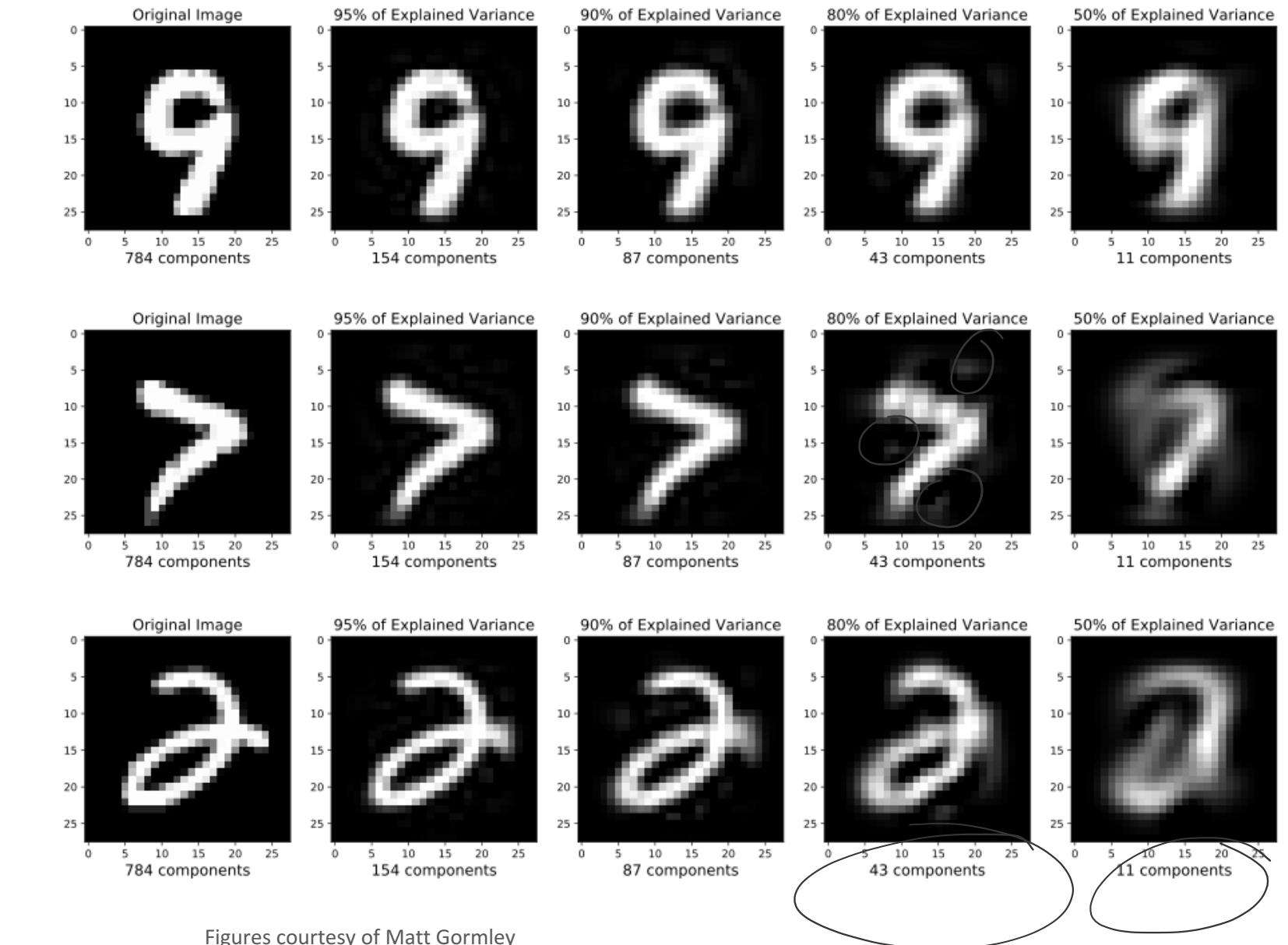
Choosing the number of PCs

- Define a percentage of explained variance for the i^{th} PC:

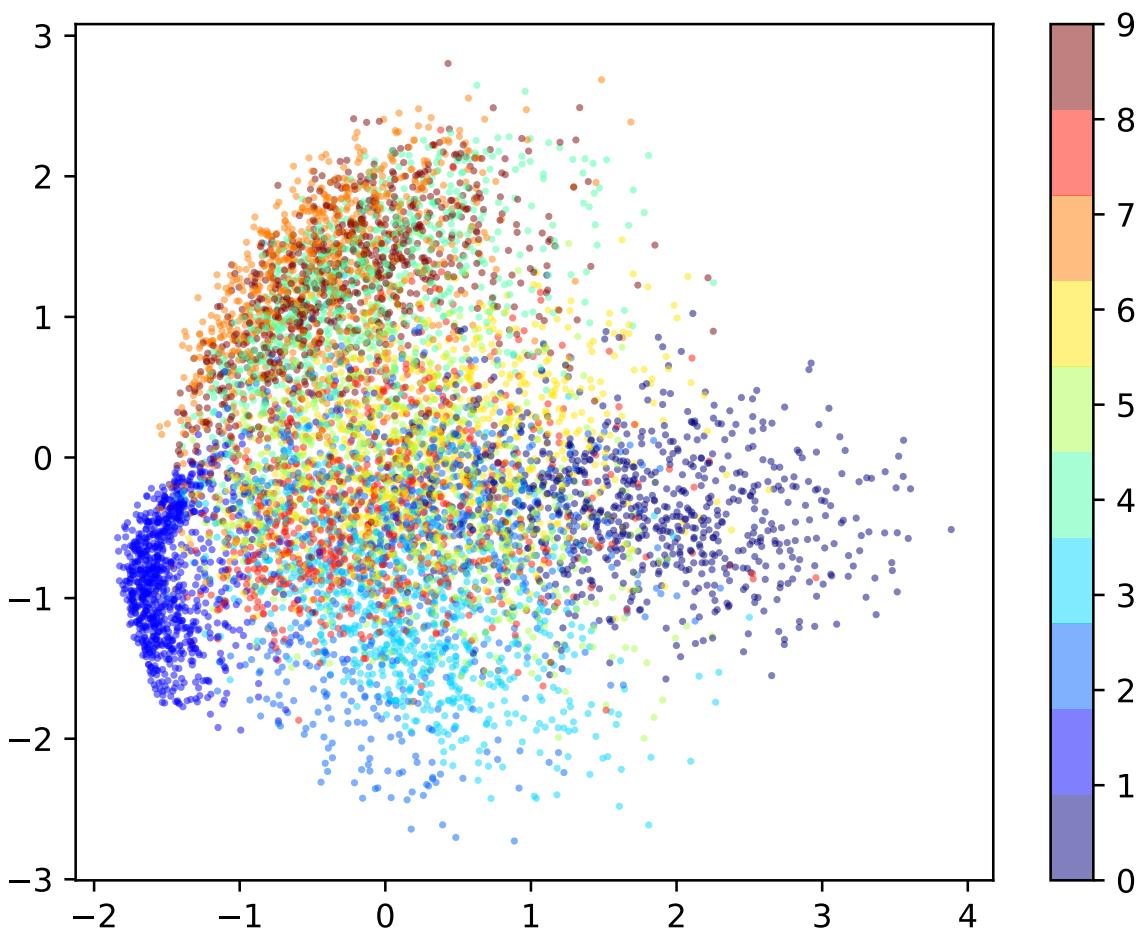
$$\lambda_i / \sum \lambda_j$$

- Select all PCs above some threshold of explained variance, e.g., 5%
- Keep selecting PCs until the total explained variance exceeds some threshold, e.g., 90%
- Evaluate on some downstream metric

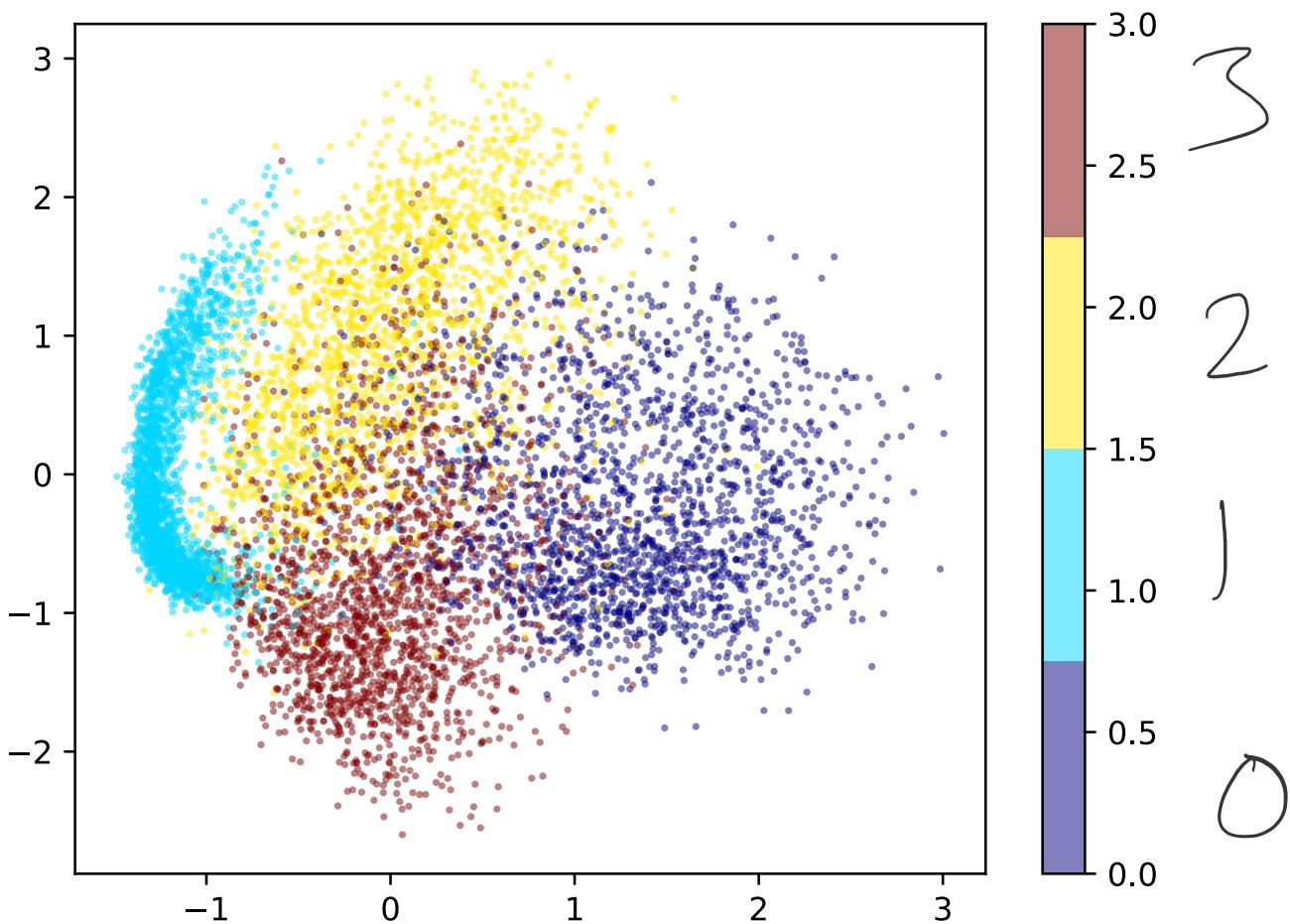
PCA Example: MNIST Digits



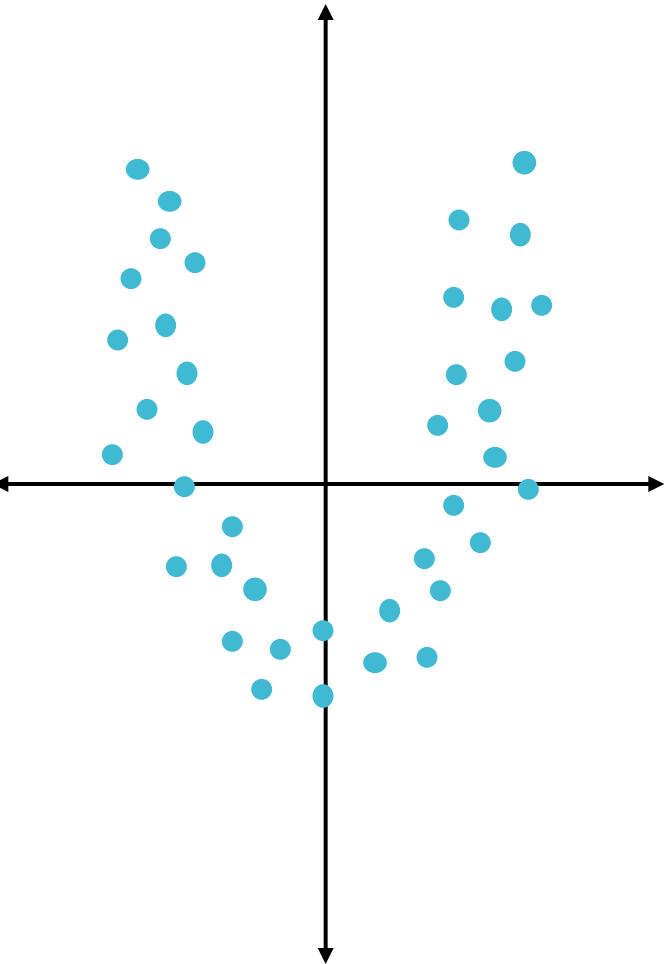
PCA Example: MNIST Digits



PCA Example: MNIST Digits



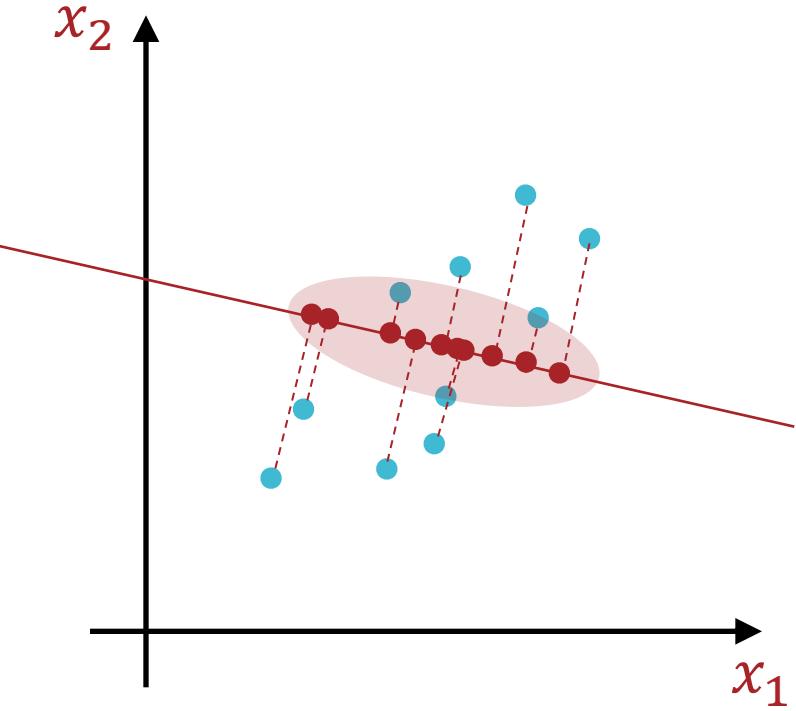
Shortcomings of PCA



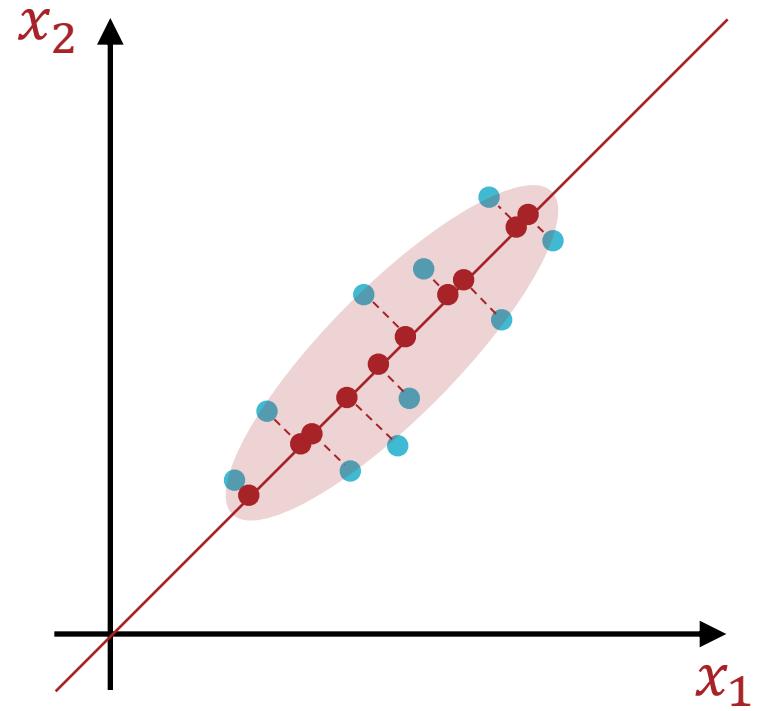
- Sometimes we don't care about variance
- PCA only finds linear combinations of our features
- Computationally intractable when k is very large
- interpretability
- orthogonality constraint

PCA Algorithm: Computational Cost

- Input: $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^n, r$
1. Center the data
 - A. Optionally, normalize the data by features so that all features are of the same scale
 2. Compute the covariance matrix $C_X = X^T X$ ($O(nk^2)$)
 3. Collect the top r eigenvectors (corresponding to the r largest eigenvalues), $P \in \mathbb{R}^{k \times r}$ ($O(k^3)$)
 4. Project the data into the space defined by P , $Z = XP$ ($O(nkr)$)



Option A



Option B

Maybe Option A isn't so bad?

Random Projections?

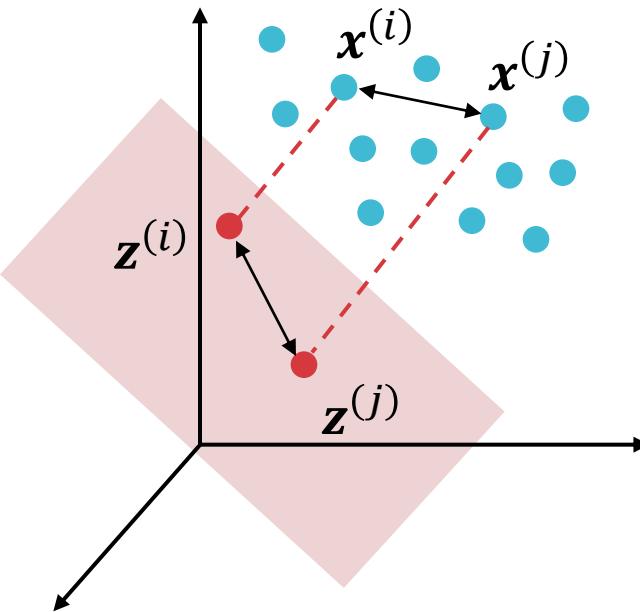
- Issue: when k is very large, computing principal components can be intractable
 - PCA also requires the entire dataset to be available, which might not be possible (e.g., online learning)
- Idea: instead of rigorously minimizing the reconstruction error/maximizing the projections' variance, just pick a random set of vectors to project onto!
 - Each element of P is sampled from a standard Gaussian distribution, $P_{i,j} \sim N(0,1)$
 - Projections are still given by $Z = XP$
- Inquiry: is this a good idea? And if so, how good?

Random Projections: Computational Cost

- Input: $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^n, r$
 1. ~~Center the data~~
~~A. Optionally, normalize the data by features so that all features are of the same scale~~
 2. ~~Compute the covariance matrix $C_x = X^T X$ ($O(nk^2)$)~~
 3. Generate P by sampling each element from a standard Gaussian distribution ($O(kr)$)
 4. Project the data into the space defined by P , $Z = XP$ ($O(nkr)$)

Preserving Relative Distances

- Reasonable desideratum: distances between points are similar before and after being projected

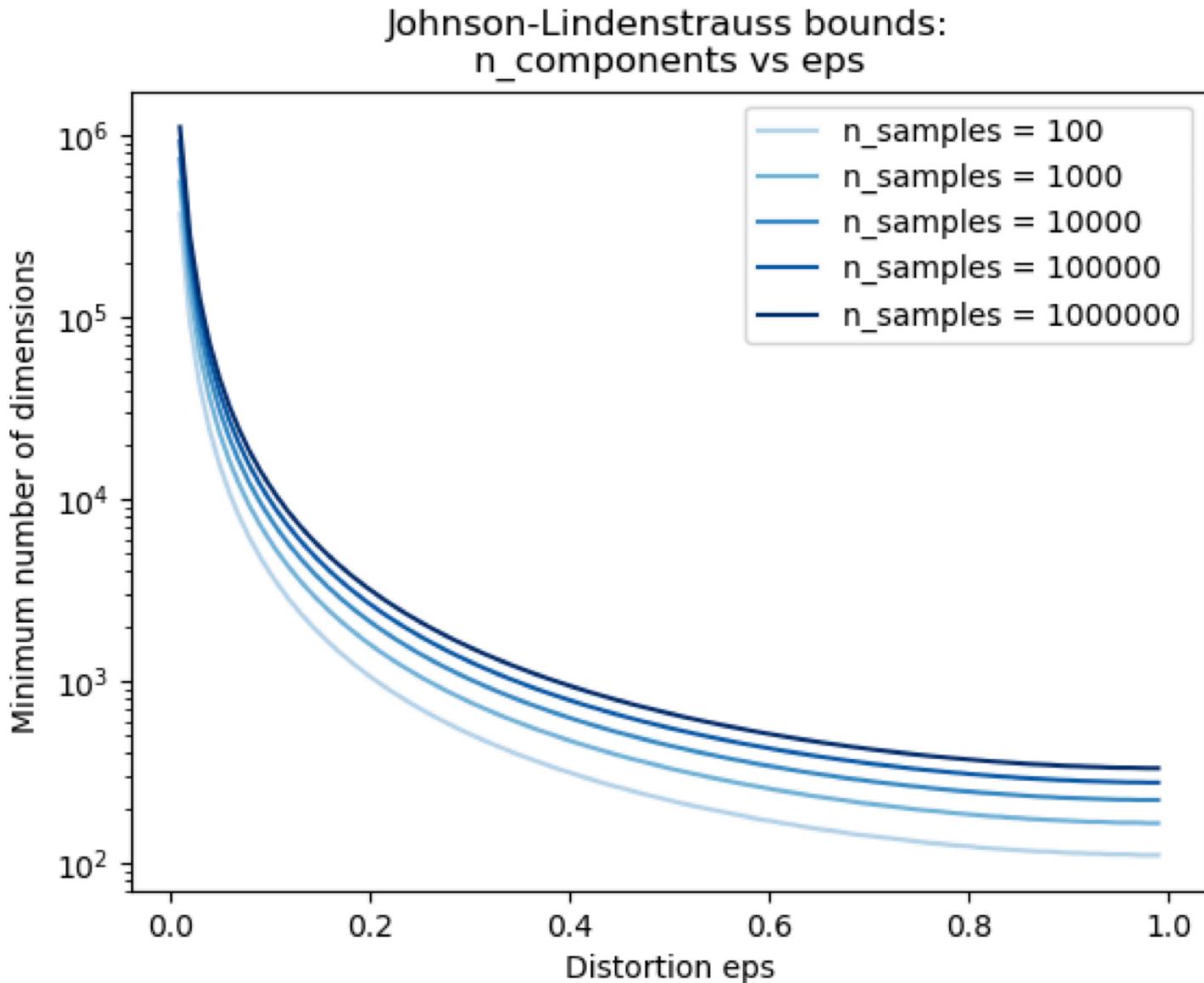


- Formally:

$$\|\mathbf{z}^{(i)} - \mathbf{z}^{(j)}\|_2^2 \in [(1 - \epsilon) \|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|_2^2, (1 + \epsilon) \|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|_2^2] \\ \forall \mathbf{x}^{(i)}, \mathbf{x}^{(j)} \in \mathcal{D}$$

Johnson- Lindenstrauss Lemma

Johnson-Lindenstrauss Lemma



Johnson- Lindenstrauss Lemma: Proof (Sketch)

Key Takeaways

- Visualization is a key component of data pre-processing
 - Visualizing big data presents unique challenges
- PCA is dimensionality reduction technique that finds an orthonormal latent representation
 - Minimizes reconstruction error \leftrightarrow maximizing the projection variance
 - Computationally expensive (cubic in the number of features)
- Random projections are an efficient alternative
 - Johnson-Lindenstrauss lemma provides theoretical guarantee for preserving distances

Key Takeaways

- Visualization is a key component of data pre-processing
 - Visualizing big data presents unique challenges
- PCA is dimensionality reduction technique that finds an orthonormal latent representation
 - Minimizes reconstruction error \leftrightarrow maximizing the projection variance
 - Constrained to find projections that are linear combinations of the existing features
 - Computationally expensive (cubic in the number of features)