

- Week 4
- **Neuroscience Introduction**
 - The brain
 - everything that happens between Stimuli and Behavior
 - Numbers of neurons
 - Studying the brain in humans
 - fMRI scanner
 - changes in blood oxygen
 - human brain
 - Relating neuronal responses to properties of an animal and its environment
 - Fine-scale sensory tuning
- **Exploratory Data Analysis**
- Model
 - Raw data
 - extracted signals
 - analysis
 - visualization
 - sharing
 - exploring
 - Interactive feedback
- Methods
 - Supervised methods
 - predict our data as a function of other data
 - Unsupervised methods
 - find structure in the data on its own
- Time series
 - supervised (regression and tuning)
 - time series as a function of some other variable
 - unsupervised (dimensionality reduction and clustering)
 - identify simple and more compact representations
 - aid our understanding
- Clustering for preprocessing
 - raw data is complex and high-dimensional
 - clustering finds collections of inputs that are similar to one another
 - these groups of clusters may be the more meaningful “unit” of measurement
- Clustering to find waveforms associated with individual neurons based on their traces across multiple electrodes
- Dimensionality reduction
 - Raw data is complex and high-dimensional
 - Dimensionality reduction describes the data using a simpler, more compact representation
 - This representation may make interesting patterns in the data more clear or easier to see
- **PCA Overview**
- 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, e.g., feature hashing
- 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
 - PCA - reconstruct 2D data via 2D data with single degree of freedom. Evaluate reconstructions (represented by blue line) by Euclidean distances
 - Linear Regression - predict y from x. Evaluate accuracy of predictions (represented by blue line) by vertical distances between points and the line
- 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 Assumptions and Solution
- PCA Formulation
 - PCA: find lower-dimensional representation of raw data
 - X is $n \times d$ (raw data)
 - $Z = XP$ is $n \times k$ (reduced representation, PCA 'scores')
 - P is $d \times k$ (columns are k principal components)
 - Linearity assumption ($Z = XP$) simplifies problem
 - Variance constraints
- Given training points
 - X matrix storing points
 - $x_j(i)$: j th feature for i 'th point
 - mau_j : mean of j 'th feature
 - Variance of 1st feature
 - the sum of the squared difference between each sample and the mean, divided by the number of samples
 - Variance of 1st feature (assuming zero mean)
 - drop the mau term since by assumption it equals zero
 - Covariance of 1st and 2nd features (assuming zero mean)
 - computing the product of the two feature values for each data point, taking the sum of these products, and finally dividing by n , the number of samples that we have
 - large positive covariance indicates that the two feature are highly correlated
 - large negative covariance indicates that the two features are highly anti correlated
 - Covariance Matrix
 - Covariance matrix generalizes this idea for many features
 - $C_x = (1/n)X'X$
 - $d \times d$ matrix

- each entry stores pairwise covariance information about the d features
- 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
 - X is $n \times d$ (raw data)
 - $Z = XP$ is $n \times k$ (reduced representation, PCA 'scores')
 - 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 C_z are zero
 - Rank-ordered features by variance, i.e., sorted diagonals of C_z
 - Variance of the first feature in the reduced dimension is the largest, followed by the variance in the second feature, and so on
- PCA Solution
 - All covariance matrices have an eigendecomposition
 - $C_x = UAU'$ (eigendecomposition)
 - U is $d \times d$ (columns are eigenvectors, sorted by their eigenvalues)
 - A is $d \times d$ (diagonals are eigenvalues, off-diagonal 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 λ_1)
- Choosing k
 - How should we pick the dimension of the new representation?
 - Visualization: Pick top 2 or top 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
 - can choose k such that we retain fraction of the variance, e.g.
- 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
- PCA Algorithm
- Orthogonal and Orthonormal Vectors
 - Orthogonal vectors are perpendicular to each other
 - equivalently, their dot product equals zero
 - Orthonormal vectors are orthogonal and have unit norm
 - a and b are orthonormal, but b and 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 $i = \{1, \dots, k\}$

- find direction of max variance that is orthonormal to previously selected directions, project onto this direction
- locations along this direction are the feature in new representation
- PCA Derivation
- Eigendecomposition
- All covariance matrices have an eigendecomposition
 - $C_x = U\Lambda U^T$ (eigendecomposition)
 - U is $d \times d$ (columns are eigenvectors, sorted by their eigenvalues)
 - Λ is $d \times d$ (diagonals are eigenvalues, off-diagonals are zero)
- Eigenvector/Eigenvalue equation: $C_x u = \lambda u$
 - By definition $u^T u = 1$ (unit norm)
- PCA Formulation
 - PCA: find lower-dimensional representation of raw data
 - X is $n \times d$ (raw data)
 - $Z = XP$ is $n \times k$ (reduced representation, PCA 'scores')
 - P is $d \times k$ (columns are k principal components)
 - Variance / Covariance constraints
- PCA Formulation, $k = 1$
 - PCA: find one-dimensional representation of raw data
 - X is $n \times d$ (raw data)
 - $z = Xp$ is $n \times 1$ (reduced representation, PCA 'scores')
 - p is $d \times 1$ (columns are k principal components)
 - variance = $\frac{1}{n} \sum z^2$ = square Euclidean norm of the vector z / n
 - Goal: maximizes variance
 - resulting principal component p to be unit norm
 - noting the relationship between the Euclidean distance and the dot product, we can rewrite the variance as a dot product $z^T z$
 - relationship between Euclidean distance and dot product = $(1/n)z^T z$

Goal: Maximizes variance, i.e., $\max_p \sigma_z^2$ where $\|p\|_2 = 1$

	$\sigma_z^2 = \frac{1}{n} \ z\ _2^2$
Relationship between Euclidean distance and dot product	$= \frac{1}{n} z^T z$
Definition: $z = Xp$	$= \frac{1}{n} (Xp)^T (Xp)$
Transpose property: $(Xp)^T = p^T X^T$; associativity of multiply	$= \frac{1}{n} p^T X^T X p$
Definition: $C_x = \frac{1}{n} X^T X$	$= p^T C_x p$

Restated Goal: $\max_p p^T C_x p$ where $\|p\|_2 = 1$

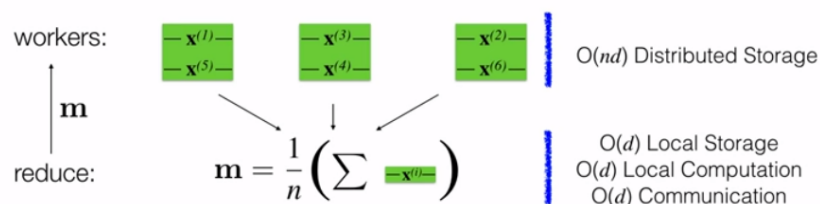
- Connection to Eigenvectors
 - Recall eigenvector/eigenvalue equation: $C_x u = \lambda u$
 - By definition $u^T u = 1$, and thus $u^T C_x u = \lambda$
 - But this is the expression we're optimizing, and thus maximal variance achieved when p is top eigenvector of C_x
- Distributed PCA

- Computing PCA Solution
 - Given: $n \times d$ matrix of uncentered raw data
 - Compute $k \ll d$ dimensional representation
 - Step 1: Center Data
 - computing the mean of each feature and subtracting the mean
 - Step 2: Compute Covariance or Scatter Matrix
 - $(1/n) * X' * X$ versus $X' * X$
 - Step 3: Eigendecomposition
 - Step 4: Compute PCA Scores
- PCA at Scale
 - Case 1 Big n and small d
 - $O(d^2)$ local storage, $O(d^3)$ local computation, $O(dk)$ communication
 - Similar strategy as closed-form linear regression
 - Case 2 Big n and Big d
 - $O(d)$ local storage and computation on workers, $O(dk)$ communication
 - Iterative algorithm
- Step 0: Data Parallel storage
 - Example: $n = 6$; 3 workers
 - workers: $O(nd)$ Distributed Storage
- Step 1: Center Data
 - Compute d feature means, m element of \mathbb{R}^d
 - Example $n = 6$; 3 workers
 - workers
 - reduce
 - $m = (1/n) * (\text{summation features})$

Step 1: Center Data

- Compute d feature means, $\mathbf{m} \in \mathbb{R}^d$
- Communicate \mathbf{m} to all workers

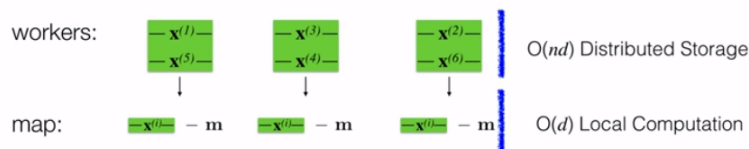
Example: $n = 6$; 3 workers



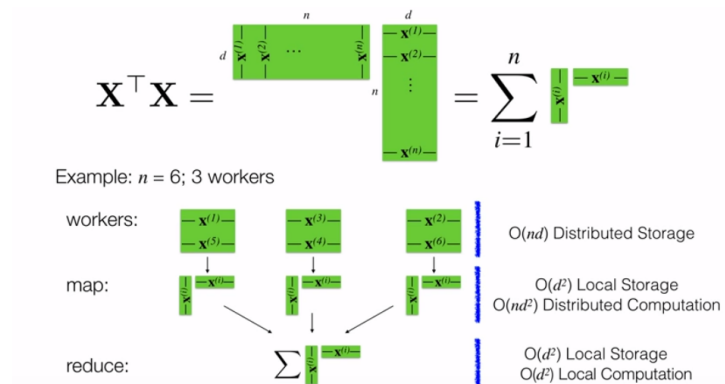
Step 1: Center Data

- Compute d feature means, $\mathbf{m} \in \mathbb{R}^d$
- Communicate \mathbf{m} to all workers
- Subtract \mathbf{m} from each data point

Example: $n = 6$; 3 workers



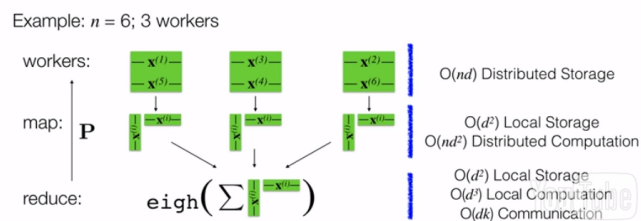
- Step 2: Compute Scatter Matrix ($\mathbf{X}^T \mathbf{X}$)
 - Compute matrix product via outer products (just like we did for closed-form linear regression!)



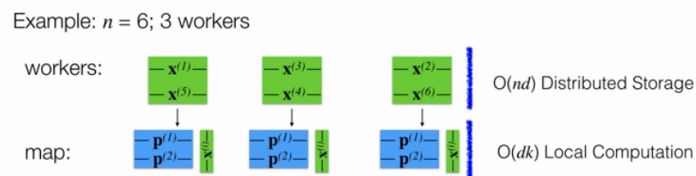
- Step 3: Eigendecomposition
 - perform locally since d is small

Step 3: Eigendecomposition

- Perform locally since d is small
- Communicate k principal components ($\mathbf{P} \in \mathbb{R}^{d \times k}$) to workers



- Step 4: Compute PCA Scores
 - Multiply each point by principal components, \mathbf{P}
 - workers \rightarrow map



- Distributed PCA, Part 2 (Optional)
- PCA at Scale
 - Case 1: Big n and Small d
 - $O(d^2)$ local storage, $O(d^2)$ local computation, $O(dk)$ communication
 - Similar strategy as closed-form linear regression
 - Case 2: Big n and Big d
 - $O(d)$ local storage and computation on workers, $O(dk)$ communication

- Iterative algorithm
- An Iterative Approach
 - We can use algorithms that rely on a sequence of matrix-vector products to compute top k eigenvectors (P)
 - E.g., Krylov subspace or random projection methods
 - Krylov subspace methods (used in MLlib) iteratively compute $X'Xv$ for some v element \mathbb{R}^d provided by the method
 - Requires $O(k)$ passes over data, $O(d)$ local storage on workers
 - We don't need to compute the covariance matrix!
 - Step1: communicate v_i element \mathbb{R}^d to all workers
 - Step2: Compute $q_i = X'Xv_i$ in a distributed fashion
 - Step1: $b_{ij} = Xv_i$
 - Step2: $q_i = X'b_i$
 - Perform in single map-reduce

Repeat for $O(k)$ iterations:

1. Communicate $v_i \in \mathbb{R}^d$ to all workers
2. Compute $q_i = X^T X v_i$ in a distributed fashion
 - Step 1: $b_i = X v_i$
 - Step 2: $q_i = X^T b_i$
 - Perform in single map-reduce!
3. Driver uses q_i to update estimate of P

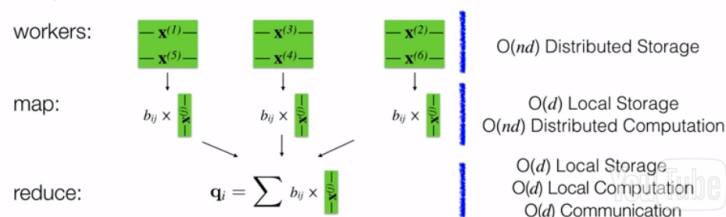
- $b_{ij} = v_i^T x^{(j)}$: each component is dot product

- q_i is a sum of rescaled data points, i.e., $q_i = \sum_{j=1}^n b_{ij} x^{(j)}$

Compute $q_i = X^T X v_i$ in a distributed fashion

- $b_{ij} = v_i^T x^{(j)}$ and $q_i = \sum_{j=1}^n b_{ij} x^{(j)}$
- Locally compute each dot product and rescale each point before summing all rescaled points in reduce step!

Example: $n = 6$; 3 workers



- `q = trainData.map(rescaleByBi).reduce(sumVectors)`