

Chemical Engineering 4H03

Mathematical Derivation of PCA

Jake Nease McMaster University

Portions of this work are copyright of ConnectMV

Objectives for this Class

- We want to compute the values p_k in each loadings vector for our principle components
 - This will complete our understanding of how PCA works
 - Will allow us to code our own PCA functions... Both in this course and)hopefully) your careers!

How?

- 1. Discussing model residuals (graphically)
- 2. Discussing model residuals (numerically)
- 3. Connection of these concepts into **Eigenvalue Decomposition** (fundamental math basis of PCA)





Error Analysis

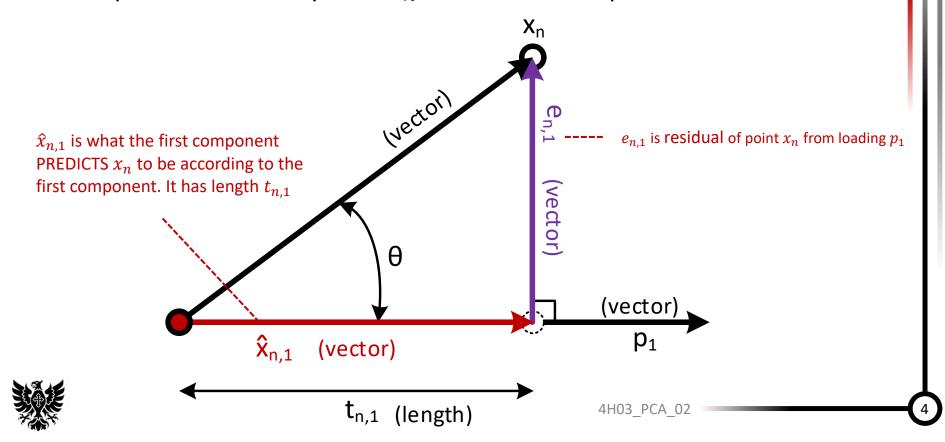
Learn from your mistakes

:(

Your PC ran into a problem and needs to restart. We're just collecting some error info, and then we'll restart for you. (0% complete)

Consider the PCA Model Prediction

- A point x_n may be reconstructed as the sum of two vectors
 - One is the vector (projection) of length t along the loading p
 - The second represents the orthogonal distance (we'll call that the residual e) from the loading to the point
- Example: residual of point x_n from first component

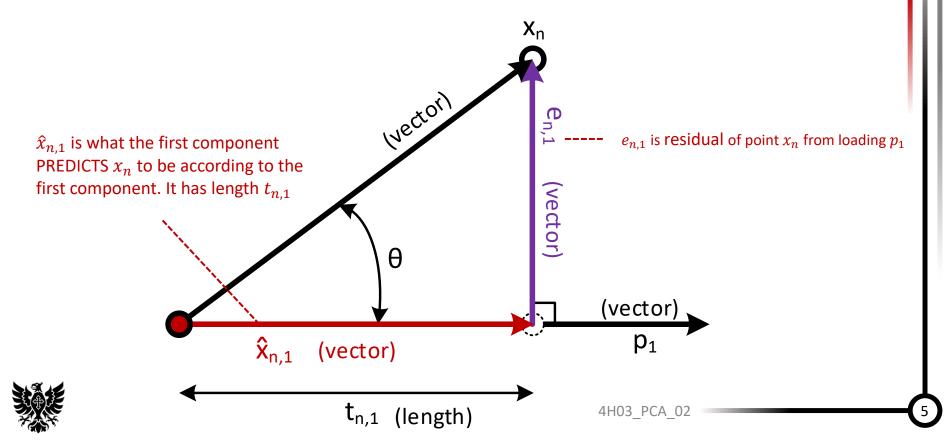


Consider the PCA Model Prediction

• We may therefore compute $\hat{x}_{n,1}$ as:

$$\hat{x}_{n,1}^T = t_{n,1} \ p_1^T - \cdots - Conveniently, t_{n,1} \text{ is a scalar that represents the DISTANCE along } p_1 \dots \\ (1 \times K) = (1 \times 1)(1 \times K)$$

The value $\hat{x}_{n,1}$ is our "best" prediction of x_n according to the first component

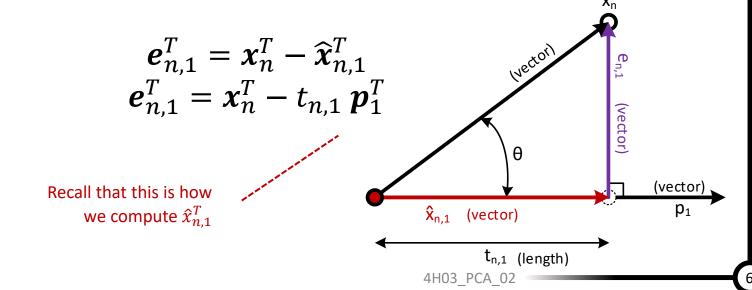


Calculating the Residual

• We can therefore compute x_n as the sum of the prediction of x_n according to the first component plus the **residual error** left behind from that component:

$$\boldsymbol{x}_n = \widehat{\boldsymbol{x}}_{n,1} + \boldsymbol{e}_{n,1}$$

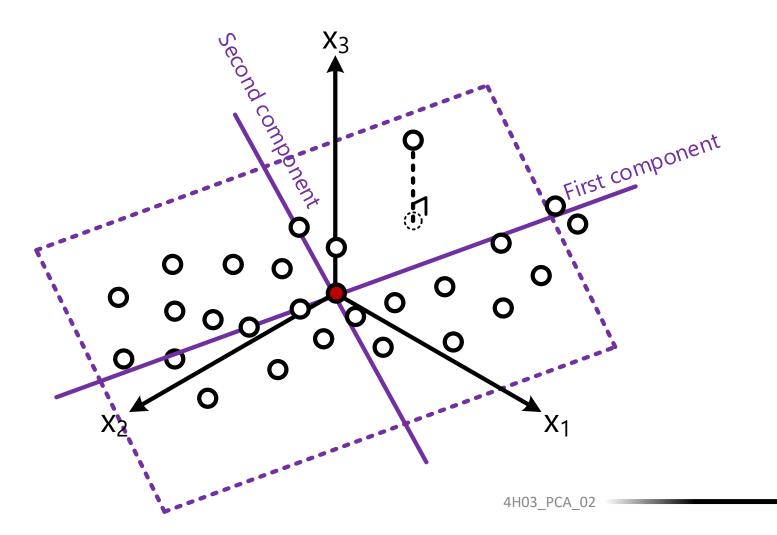
• So, the entire **ROW** of residuals for a point x_n is $e_{n,1}$





Question

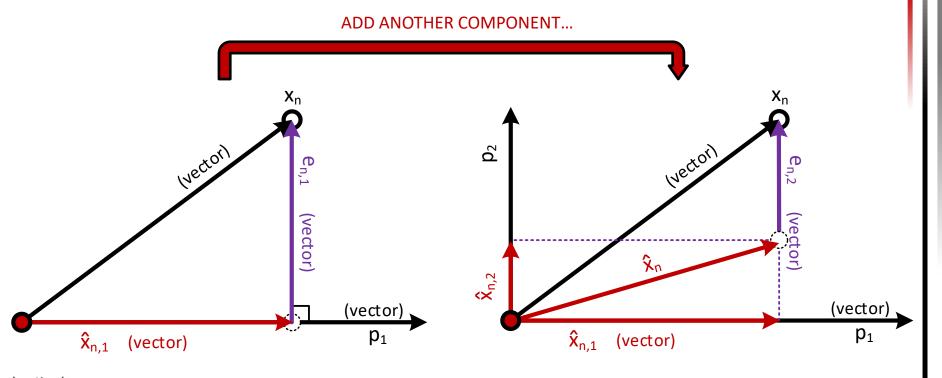
• So if that is how the residual is computed... what do you expect to happen if we add another component?





Question

- So if that is how the residual is computed... what do you expect to happen if we add another component?
 - We should find the residual error $e_{n,2}$ is LOWER since we have another dimension attempting to explain the data:





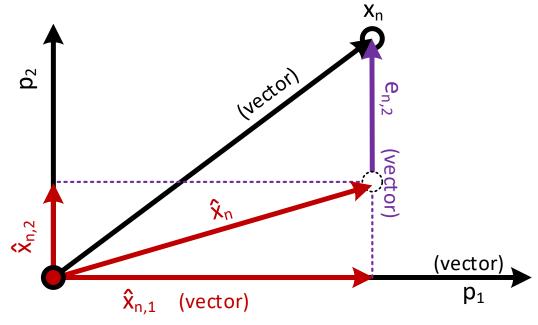
Overall Prediction of \hat{X}

• At the end of the day, if we have A components, the model predicts \hat{X} as:

$$\hat{X} = TP^T$$

MUCH easier than computing each individual \hat{x} and adding them up (math is the same)

$$(N \times K) = (N \times A) \times (A \times K)$$





Quantifying Residuals

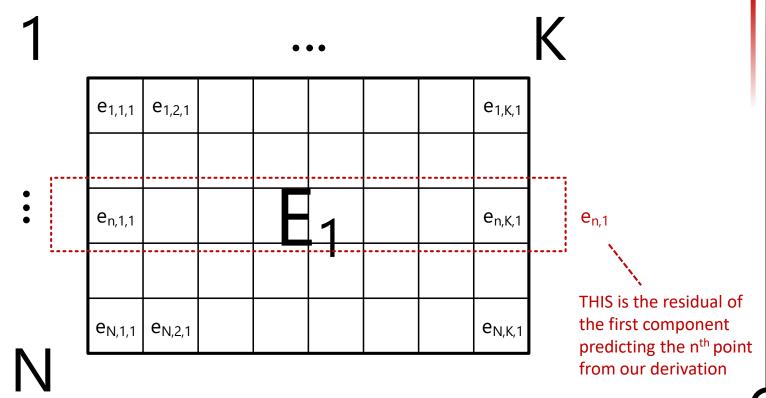
- Recall each residual $e_{n,1}$ is a row from ${f X}$
- We can assemble all residuals into a matrix that is the same size as X and represents the residual of our loading to each variable
 - For our first component, we'll call this $E \in \mathbb{R}^{N \times K}$

1	• • •								K
	e _{1,1,1}	e _{1,2,1}						e _{1,K,1}	
•	e _{n,1,1}				1			e _{n,K,1}	
.	e _{N,1,1}	e _{N,2,1}						e _{N,K,1}	



Quantifying Residuals

- Recall each residual $e_{n,1}$ is a row from ${\bf X}$
- We can assemble all residuals into a matrix that is the same size as X and represents the residual of our loading to each variable
 - For our first component, we'll call this $E \in \mathbb{R}^{N \times K}$





So What Do We Know?

• We know that we can compute \widehat{X} (all predicted measurements) using the convenient expression:

$$\hat{X} = TP^T$$

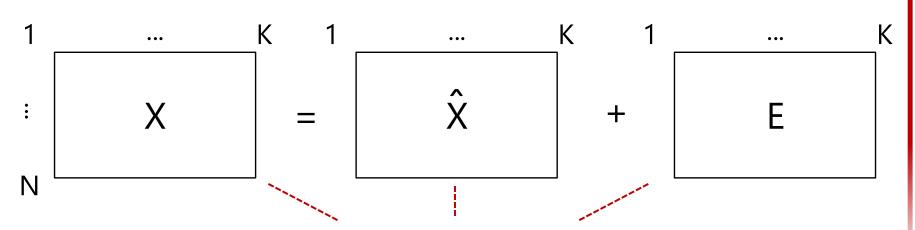
• We know that the TOTAL error of our model predictions \hat{X} are the difference between X and the model prediction:

$$X = \hat{X} + E \implies X = TP^T + E \implies E = X - TP^T$$

• We also know that the more components (columns in T) we use, the close \widehat{X} becomes to X, and thus the smaller E becomes!



So What Do We Know?



It is important that these are all the same dimensions, else we could not add them!



$$(N \times K) = (N \times A) \times (A \times K)$$

Looking Ahead...

- So, if I were to ask you right now how we calculate the values in matrix *P*...
 - AKA the loadings, for component a: \boldsymbol{p}_a
- Where do you think this is going...?
 - Yup, you'd best believe that we want to minimize the variance in E. That is... minimize the sum of squared-errors*
 - If we can find the values of P that make E as small as possible (for some known number of components A), we can say that the LV model as described by P optimally explains the variance in X

In other words: perform a regression.





Determining Model Accuracy

If the glove fits...



Determining Model Fit

- You may ask why we are bothering to compute \widehat{X} since we already have the data for X
 - Good question… I have two answers:
 - 1. To BACK-CALCULATE values of X from a given or desired score T (reverse-engineering approach)
 - 2. To compute the **goodness of fit** of our model!
- Recall the metric R^2 ?
 - $-R_{a=0}^2 = 0$ (no variance explained)
 - $-R_a^2$ ALWAYS \uparrow as $a \uparrow$ (less each time)

$$- R_{a=1}^2 < R_{a=2}^2 < \dots < R_{a=A}^2$$

− Overall, $R^2 \le 1$

components are added
$$R_a^2 = 1 - \frac{\mathcal{V}(X - \hat{X})}{22(X)}$$

 E_a are residuals after \hat{X} is predicted using a components

 R^2 of the model after a

$$R_a^2 = 1 - \frac{\mathcal{V}(E_a)}{\mathcal{V}(X)}$$



Workshop – Model Fit of Pastries

- Now that we know how to do it...
 - Use PCA function to determine loadings/scores for pastries
 - Compute R^2 using different numbers of components
 - Live demo in MATLAB

$$R_a^2 = 1 - \frac{\mathcal{V}(X - \widehat{X})}{\mathcal{V}(X)}$$

$$R_a^2 = 1 - \frac{\mathcal{V}(E_a)}{\mathcal{V}(X)}$$



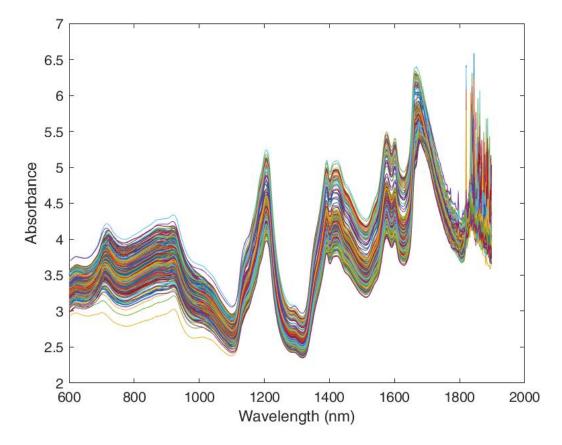
Example: Spectral Data

- Recall our spectral data set (N = 460, K = 650)
 - It turns out we need ONLY THREE latent variables to describe a significant portion of the variance in the data

•
$$R_{a=1}^2 = 0.737$$

•
$$R_{a=2}^2 = 0.922$$

- Adds 18.5%
- $R_{a=3}^2 = 0.942$
 - Adds 2%





Column Residuals

- R² may be calculated for each **column** instead of the entire data set
 - Represents the predictive capability of our LVM space for one variable in particular

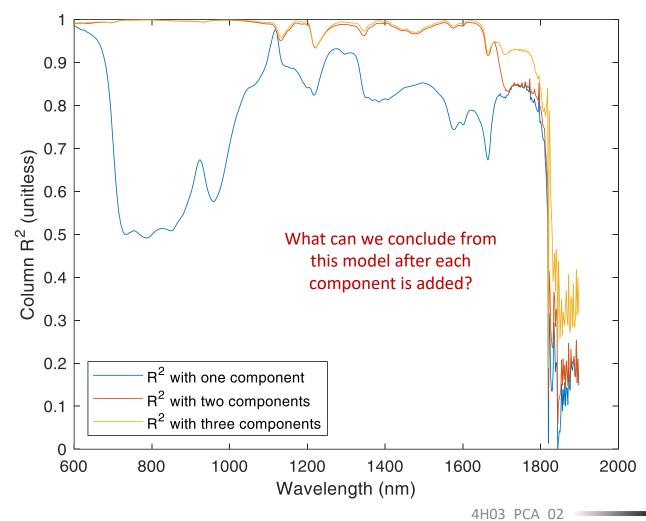
•
$$R_k^2 = 1 - \frac{\mathcal{V}(x_k - \widehat{x}_k)}{\mathcal{V}(x_k)}$$

- More components will always increase how well model explains the column
- $-R_k^2 = 0$ when there are no components (no shocker)
- $-R_k^2$ increases for every component added
- **Question**: does it increase a lot if the new component has a loading $p_{k,a} \approx 0$?



Column Residuals

Example for spectral data in MATLAB



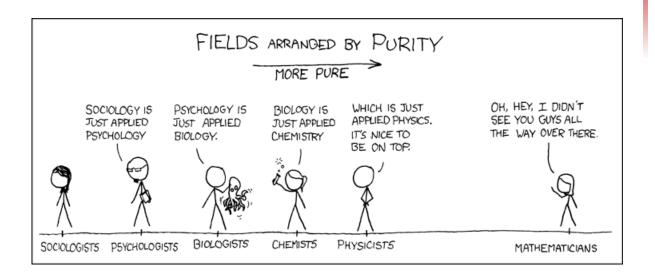


(20)



Derivation of PCA Components

It's math... IN (eigen) SPAAAAAAAAACE!!



Optimization is Everywhere

- Similar to what we did in the regression section, we can set up an **optimization** to find \boldsymbol{p}_a
 - Let's just do this for the first component for now, so p_1
- What does this optimization look like?

IN GENERAL:

$$\min_{???} \phi = f(???)$$

- For PCA, what are some possible...
 - Objective functions?
 - Constraints?



- **FIRST**: WHY?
 - Because it gives us great insight as to what we are doing
 - Because it relates closely to regression
- SECOND: HOW?
 - On the board for now
 - Summary of critical results follow, but...
 - Will make no sense unless you can follow the math ☺
- Warning: I am going to make some hand-wavy claims
 - It is not the purpose of this class to get too detailed
 - If you have questions, I'm always game ☺



$$\max_{\boldsymbol{p}_1} \boldsymbol{\phi} = \boldsymbol{t}_1^T \ \boldsymbol{t}_1$$

 $\max \phi = (X \boldsymbol{p}_1)^T X \boldsymbol{p}_1 \mid \max \phi = \boldsymbol{p}_1^T X^T X \boldsymbol{p}_1$

(subject to)

(subject to)

(subject to)

$$\boldsymbol{p}_1^T \, \boldsymbol{p}_1 = 1$$
$$\boldsymbol{t}_1 = X \, \boldsymbol{p}_1$$

$$\boldsymbol{p}_1^T \, \boldsymbol{p}_1 = 1$$

$$\boldsymbol{p}_1^T \, \boldsymbol{p}_1 = 1$$

- This is **nearly** a convex optimization program
 - You'll notice that the max looks like a sum of squares
 - This is very similar to our regression
- The constraint really screws us up
 - However, the fact that it is an equality constraint gives us some hope... Maybe there is a way to reflect it in the objective?



- Our strategy is to ELEVATE the constraint $p_1^T p_1 = 1$ to the objective and penalize a deviation from that unity
 - However, instead of choosing some arbitrary penalty factor, we will add a decision variable as the penalty factor
 - There are *reasons* for this, which we can discuss
 - The penalty factor is known as a **Lagrange Multiplier** $\lambda_1 \geq 0$

You might ask... Why not make $\lambda_1=0$? That would be equivalent to ignoring the constraint, which means p_1 can be infinite magnitude (no unique solution!

So if we deviate from $p_1^T p_1 = 1$, we *subtract* from the objective function

$$\max_{\boldsymbol{p}_1} \phi = \boldsymbol{p}_1^T X^T X \, \boldsymbol{p}_1 - \lambda_1 (\boldsymbol{p}_1^T \, \boldsymbol{p}_1 - 1)$$

(Unconstrained)



- This problem is unconstrained, and allows us to identify a stationary point (like SSE for regression)
 - That is, we take the derivative of ϕ with respect to all decision variables and set them equal to 0
 - So, we take the derivative WRT all entries in p_1 AND WRT λ_1
 - How many equations and unknowns will I have?

$$\frac{\partial \phi}{\partial \boldsymbol{p}_1} = 2X^T X \boldsymbol{p}_1 - 2\lambda_1 \boldsymbol{p}_1 = \boldsymbol{0}$$
 ———— That's a $(K \times 1)$ vector of 0

$$\frac{\partial \phi}{\partial \lambda_1} = p_1^T \ p_1 - 1 = 0$$
 ----- That's a (1 × 1) scalar of 0



- Look what happened thanks to λ_1
 - I am still obliged to choose $p_1^T p_1 = 1$
 - BUT the gradient of ϕ WRT $m{p}_1$ includes λ_1
 - Now, the only way to get $\frac{\partial \phi}{\partial p_1} = 0$ is to choose $p_1 = 0$ OR find a magical value of λ_1 and p_1 that balances both pieces of $\frac{\partial \phi}{\partial p_1}$
 - Note I can't even choose $p_1 = 0$ because the second equation would be violated!

$$\frac{\partial \phi}{\partial \boldsymbol{p}_1} = 2X^T X \boldsymbol{p}_1 - 2\lambda_1 \boldsymbol{p}_1 = \mathbf{0} \quad \text{That's a } (K \times 1) \text{ vector of 0}$$

$$\frac{\partial \phi}{\partial \lambda_1} = p_1^T \ p_1 - 1 = 0$$
 ----- That's a (1 × 1) scalar of 0



• Focusing only on the $\frac{\partial \phi}{\partial p_1}$ term:

$$X^T X \boldsymbol{p}_1 - \lambda_1 \boldsymbol{p}_1 = \boldsymbol{0}$$
 ----- I divided out the 2. Hurray!

• Factoring out p_1 on the RHS of the matrices:

$$(X^TX - \lambda_1 I_K) \boldsymbol{p}_1 = \boldsymbol{0}$$
 ----- Also note that $\boldsymbol{p}_1^T \boldsymbol{p}_1 = 1$

- Hang on just a hot minute...
 - Where did this magical I_K come from?
 - Does that expression look familiar?



- CONCLUSION: The first component of $X(p_1)$, that...
 - Has unit length
 - Describes maximum variance in X
 - Correspond to scores t_1
- ... Is EXACTLY equal to the first Eigenvector of X^TX
 - It has corresponding Eigenvalue λ_1 such that $X^T X \mathbf{p}_1 = \lambda_1 \mathbf{p}_1$
 - $-\lambda_1$ is exactly the variance described by the first component
- Subsequent components are the remaining eigenvectors of X^TX
 - Provided $p_i \cdot p_j = 0 \ \forall \{i, j\}$ (all components are orthogonal)
 - The eigenvectors λ_i for each component \boldsymbol{p}_i are the variances in X^TX in the direction \boldsymbol{p}_i



Eigenvalue Decomposition Summary

• For the first component, our optimization problem to choose the first loadings vector p_1 is:

These are the same...

$$\max_{m{p}_1} m{\phi} = m{t}_1^T \ m{t}_1$$
 OR $\min_{m{p}_1} m{\phi} = \sum_{i=1}^N e_{i,1}^2$ (subject to) $m{p}_1^T \ m{p}_1 = 1 \ \cdots \ \sum_{p_1 \text{ to have unit length}}^{\text{Recall that we are forcing}} m{t}_1 = X m{p}_1$

• ALSO note that p_1 has K values ($p_1 \in \mathbb{R}^{K \times 1}$)



Eigenvalue Decomposition Summary

- Using this optimization problem and extending to multiple components (we'll skip for components 2+):
 - Loadings p_a are the **Eigenvectors** of X^TX
 - Scores $\boldsymbol{t}_a = X\boldsymbol{p}_a$ *OR* T = XP
 - Use scores to compute $\hat{X}_A = TP^T$
 - Determine residuals $E_A = X \hat{X}_A$
 - Finally, can compute model fit as $R^2 = 1 \frac{\mathcal{V}(E_A)}{\mathcal{V}(X)}$
- Fun facts about this decomposition...
 - Eigenvalues λ_a are the variances of the scores t_a
 - Sum of all eigenvalues of X^TX is trace(X) is $\mathcal{V}(X^TX)$



Workshop

- So we know we can calculate the loadings as the Eigenvectors of $X^TX...$
 - Explain why this works very well for "tall" matrices where $N\gg K$
 - Explain what might happen if I have a large number of columns such as the spectral data set... $N \ll K$



Eigenvalue Decomposition Remarks

- If you have a "short" matrix X, you might want to consider avoiding $X^TX \in \mathbb{R}^{K \times K}$
 - Instead, recognize that the Eigenvectors of $XX^T \in \mathbb{R}^{N \times N}$ are actually the scores of X, just not necessarily scaled correctly
- So our (more computationally efficient) strategy becomes:
 - Determine t_a as the Eigenvectors of XX^T
 - Determine "unscaled" loadings as $p_a = Xt_a$
 - Normalize loadings as $p_a = \frac{p_a}{\|p_a\|}$
 - Re-calculate scores as $t_a = X p_a$



Eigenvalue Decomposition Remarks

- There are several disadvantages of EV Decomposition to determine the loadings of a PCA model
 - We calculate all K Eigenvectors when we usually want $A \ll K$
 - Any missing data completely breaks this method
 - WORKSHOP: how would you handle missing data?

- Moreover, we have problems with large matrices
 - Calculating X^TX can be cumbersome
 - Large data sets can cause overflow/loss of memory
 - We have to keep X in memory to compute scores



Eigenvalue Decomposition Remarks

- But there are distinct advantages too
 - EV decomposition gives us terrific insight into what a PCA model means and what we are doing to extract that info
 - An amazing application of an abstract math concept
 - All properties of PCA are derived from this method
 - Are slightly more accurate than convergence-based methods such as NIPALS
 - Only source of error is numerical precision
 - Error is "spread out" over all components
 - In NIPALS, the algorithm becomes less accurate the more components that are added



Final Remarks

- OK! So we have seen some cool stuff here
 - Geometric interpretation of model residuals
 - Mathematical interpretation of model residuals
 - Eigenvalue decomposition derivation
 - Connection of EV decomposition to PCA statistics
- Where are we going next?
 - NIPALS algorithm, which...
 - Can be much more efficient than EV decomposition
 - Can handle missing data
 - PLS
 - Extension of PCA to also predict outputs
 - NIPALS will be coded by YOU for this one ©

