

Linear Regression

C_p , AIC, BIC, and Adjusted R^2

Prof. Asim Tewari
IIT Bombay

$$MSE = \frac{RSS}{n}$$

$$(MSE)_{\text{Training}} = \frac{RSS_{\text{Training}}}{n}$$

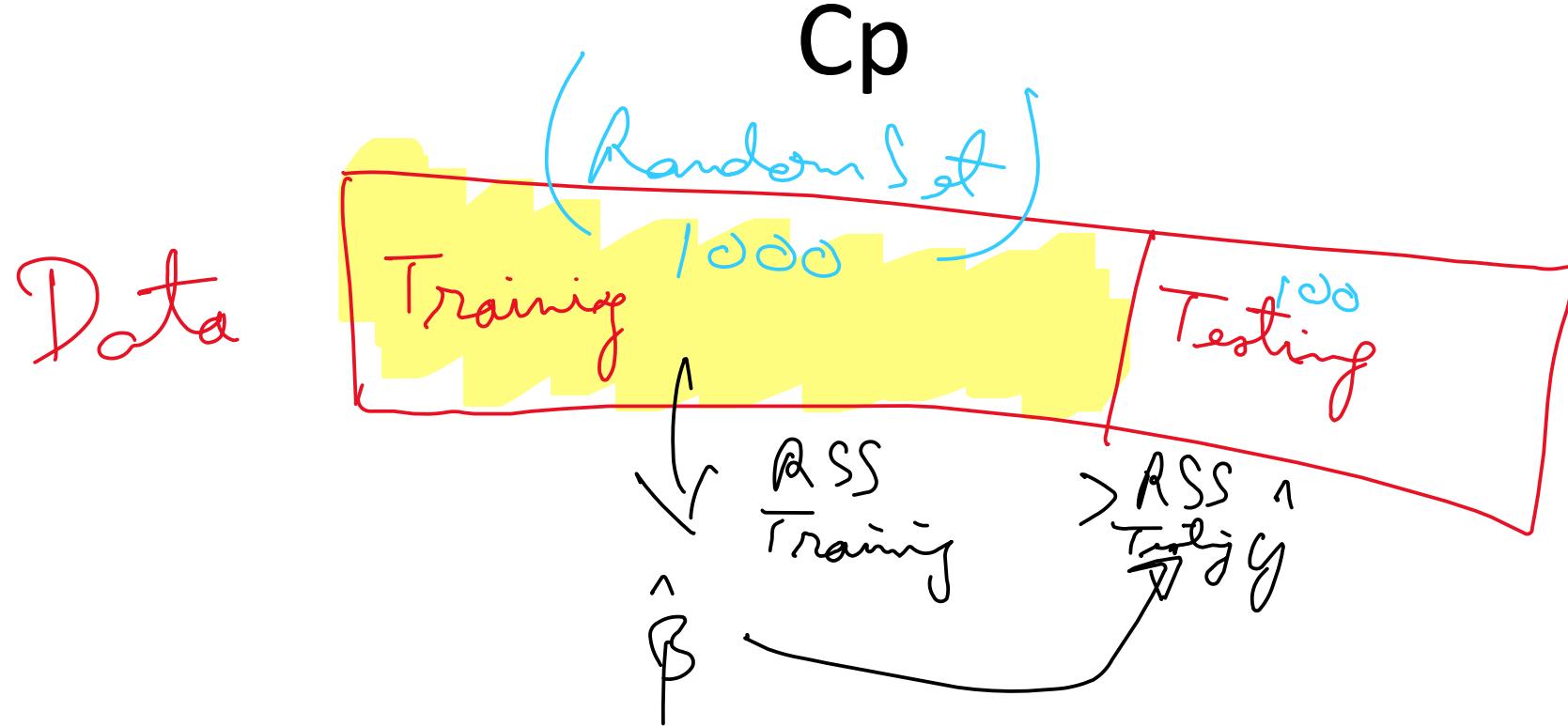
$$(MSE)_{\text{Testing}} = \frac{RSS_{\text{Testing}}}{n}$$

- For a fitted least squares model containing d predictors, the C_p estimate of test MSE is computed using the equation

$$\underline{C_p} = \frac{1}{n} (RSS + 2d\hat{\sigma}^2)$$

Estimator
 of MSE of
 Testing

- where $\hat{\sigma}^2$ is an estimate of the variance of the error associated with each response measurement



AIC

- The AIC criterion is defined for a large class of models fit by maximum likelihood. In the case of the model

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \epsilon$$

with Gaussian errors, maximum likelihood and least squares are the same thing.

$$\text{AIC} = \frac{1}{n\hat{\sigma}^2} (\text{RSS} + 2d\hat{\sigma}^2)$$

BIC

- BIC is derived from a Bayesian point of view, but ends up looking similar to C_p (and AIC) as well. For the least squares model with d predictors, the BIC is, up to irrelevant constants, given by

$$\log(n) > 2 \Rightarrow n > 7$$

$$\text{BIC} = \frac{1}{n} (\text{RSS} + \underbrace{\log(n)d\hat{\sigma}^2}_{2d\hat{\sigma}^2})$$

$$C_p = \frac{1}{n} (\text{RSS} + \underbrace{2d\hat{\sigma}^2}_{2d\hat{\sigma}^2})$$

Adjusted R^2

$$R^2 = 1 - \frac{RSS}{TSS} ; \quad TSS = \sum (y_i - \bar{y})^2$$

$$\text{Adjusted } R^2 = 1 - \frac{RSS/(n-d-1)}{TSS/(n-1)}$$

Choosing the Optimal Model

For a fitted least squares model containing d predictors

- C_p

$$\frac{\text{RSS}}{n} = \text{MSE}_{\text{of}} \text{ to sig}$$

$$C_p = \frac{1}{n} (\text{RSS} + 2d\hat{\sigma}^2) \text{ a measure for MSE for testing}$$

- Akaike information criterion (AIC)

$$\text{AIC} = \frac{1}{n\hat{\sigma}^2} (\text{RSS} + 2d\hat{\sigma}^2)$$

- Bayesian information (BIC)

$$\text{BIC} = \frac{1}{n} (\text{RSS} + \log(n)d\hat{\sigma}^2)$$

- Adjusted R^2

where $\hat{\sigma}^2$ is an estimate of the variance of the error associated with each response measurement

$$\text{Adjusted } R^2 = 1 - \frac{\text{RSS}/(n - d - 1)}{\text{TSS}/(n - 1)}$$

Regression

$$y = f(x) + \epsilon$$

↑

Find 'f'

$$\hat{y} = h(x)$$

↳ Parametric model/function

Data

$$D = \left\{ (\bar{x}_i, \bar{y}_i) \right\}_{i=1}^n$$

$$\text{E.g. } \hat{f}(x, \beta) = \beta_0 + \beta_1 x^{(1)} + \beta_2 x^{(2)}$$

ML is to use data to find β .

Minimize Cost Function $L(\beta)$

$$\hat{\beta} = \arg \min_{\beta} L(\beta)$$

Null Hypothesis

→ x_1, x_2, \dots are
true features or not?

Is $\beta_i = 0$?

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \epsilon$$

which of the X_i 's are actually important?

Obj: Choose the most imp. 'd' variable from the given 'p' variables.

- 1.) Subset Selection:
- 2.) Shrinkage Methods
- 3.) Dimension Reduction

Subset Selection

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

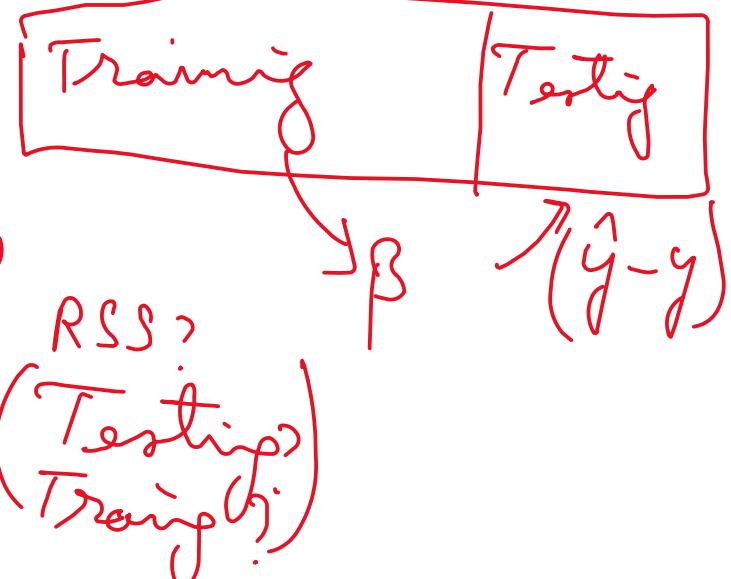
$$\left\{ \begin{array}{l} 0 \text{ variable: } \text{1 way } (y = \beta) \quad h_{\beta}(x, \beta) \\ 1 \text{ variable: } P \text{ ways } (y = \beta_0 + \beta_1 x) \\ 2 \text{ variables: } \frac{P(P-1)}{2} = {}^P C_2 \\ d \text{ variables: } {}^P C_d = \frac{\underbrace{P}_{\text{ways}}}{\underbrace{(P-d)!d!}} \\ P \text{ variables: } 1 \text{ way} \end{array} \right.$$
$$1 + {}^P C_1 + {}^P C_2 + \dots + {}^P C_d + \dots + \underbrace{{}^P C_P}_{\text{ways}} = (1+1)^P = 2^P$$

Subset Selection

Which one is better?

$$\textcircled{1.} \quad \hat{y} = \beta_0 + \beta_1 x_1$$

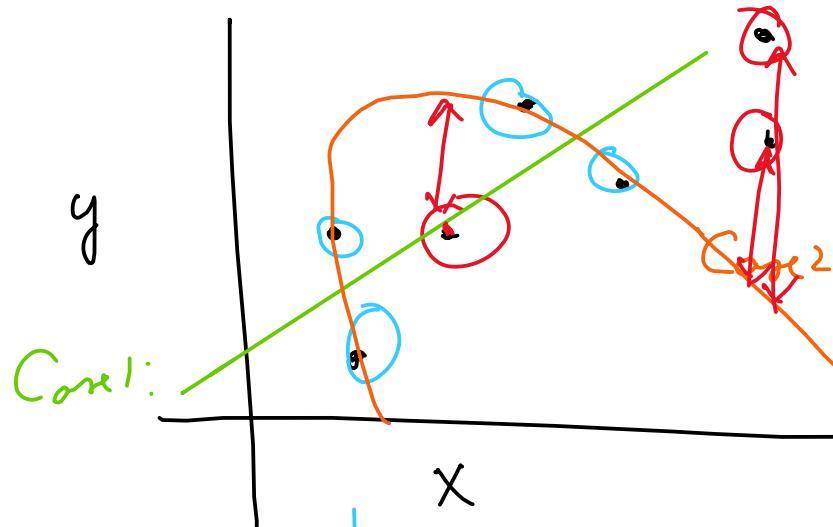
$$\textcircled{2.} \quad \hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$



$$\text{RSS of } \textcircled{1.} \geq \text{RSS of } \textcircled{2.}$$

Better \Rightarrow Lower MSE of Testing

Subset Selection



$$MSE = \frac{RSS}{n}$$

RSS of case 1 \gg RSS of case 2

$$\begin{aligned} MSE_{\text{Training}} & \quad RSS_{\text{of Testing}} \\ MSE_{\text{Testing}} & \quad < \quad RSS_{\text{of Testing}} \\ & \quad \text{of Case 1} \quad \text{of Case 2} \\ & \quad (MSE)_{\text{Test}}^{\text{Case 1}} < (MSE)_{\text{Test}}^{\text{Case 2}} \end{aligned}$$

Subset Selection

- We cannot try all the 2^p models that contain subsets of p variables.
 - Thus we have strategies for choosing the subset.
 - Various statistics can be used to judge the quality of a model. These include
 - Mallow's C p , Akaike information criterion (AIC),
 - Bayesian information criterion (BIC), and
 - Adjusted R²
- $\left. \begin{matrix} \text{MSE} \\ \text{fTest} \end{matrix} \right\}$

Subset Selection

- We cannot try all the 2^p models that contain subsets of p variables.
- Thus we have strategies for choosing the subset.
- Various statistics can be used to judge the quality of a model. These include
 - Mallow's C p , Akaike information criterion (AIC),
 - Bayesian information criterion (BIC), and
 - Adjusted R²

Strategies for choosing the subset

1. Best-Subset Selection:

Best subset regression finds for each $k \in \{0, 1, 2, \dots, p\}$ the subset of size k that gives smallest MSE for testing

2. Forward- and Backward-Stepwise Selection:

its selection starts with the intercept, and then sequentially adds into the model the predictor that most improves the fit.

Subset Selection

$k = 1, \mathcal{M}_1$
 $k = 2, \mathcal{M}_2$
 \vdots
 M_p

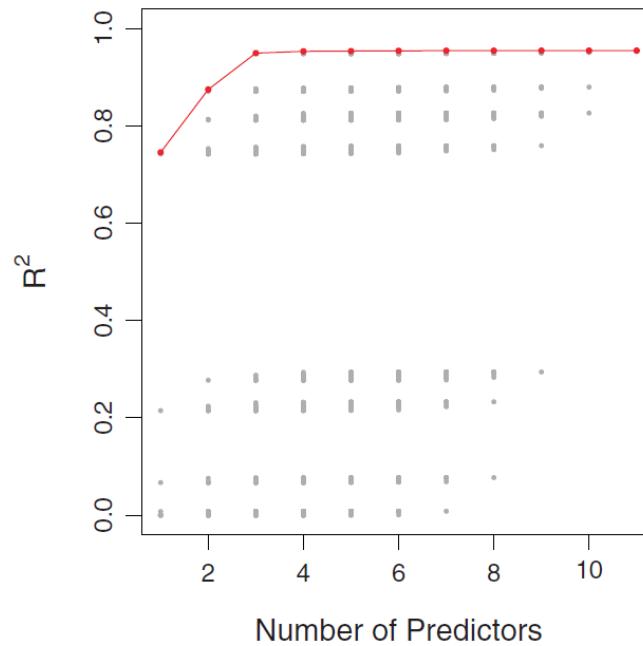
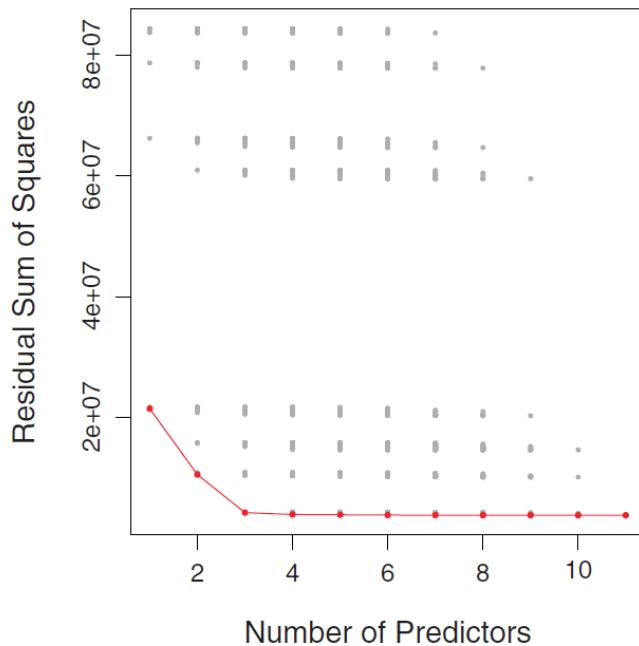
- *Best Subset Selection*

Algorithm

-
1. Let \mathcal{M}_0 denote the *null model*, which contains no predictors. This model simply predicts the sample mean for each observation.
 2. For $k = 1, 2, \dots, p$:
 - (a) Fit all $\binom{p}{k}$ models that contain exactly k predictors.
 - (b) Pick the best among these $\binom{p}{k}$ models, and call it \mathcal{M}_k . Here *best* is defined as having the smallest RSS, or equivalently largest R^2 .
 3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .

Subset Selection

- *Best Subset Selection*



For each possible model containing a subset of the ten predictors in the Credit data set, the RSS and R^2 are displayed. The red frontier tracks the best model for a given number of predictors, according to RSS and R^2 . Though the data set contains only ten predictors, the x-axis ranges from 1 to 11, since one of the variables is categorical and takes on three values, leading to the creation of two dummy variables.

Strategies for choosing the subset

1. Forward selection

- We begin with the null model—a model that contains an intercept but no predictors.
- We then fit p simple linear regressions and add to the null model the variable that results in the lowest RSS.
- We then add to that model the variable that results in the lowest RSS for the new two-variable model.
- This approach is continued until some stopping rule is satisfied.

Subset Selection

- *Stepwise Selection*
 - Forward Stepwise Selection
-

1. Let \mathcal{M}_0 denote the *null* model, which contains no predictors.
 2. For $k = 0, \dots, p - 1$:
 - (a) Consider all $p - k$ models that augment the predictors in \mathcal{M}_k with one additional predictor.
 - (b) Choose the *best* among these $p - k$ models, and call it \mathcal{M}_{k+1} . Here *best* is defined as having smallest RSS or highest R^2 .
 3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .
-

Strategies for choosing the subset

2. Backward selection

- We start with all variables in the model, and remove the variable with the largest p-value—that is, the variable
- that is the least statistically significant.
- The new $(p - 1)$ -variable model is fit, and the variable with the largest p-value is removed.
- This procedure continues until a stopping rule is reached.
- For instance, we may stop when all remaining variables have a p-value below some threshold.

Subset Selection

- *Stepwise Selection*
 - Backward Stepwise Selection
-

1. Let \mathcal{M}_p denote the *full* model, which contains all p predictors.
 2. For $k = p, p - 1, \dots, 1$:
 - (a) Consider all k models that contain all but one of the predictors in \mathcal{M}_k , for a total of $k - 1$ predictors.
 - (b) Choose the *best* among these k models, and call it \mathcal{M}_{k-1} . Here *best* is defined as having smallest RSS or highest R^2 .
 3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .
-

Strategies for choosing the subset

3. Mixed selection

- This is a combination of forward and backward selection.
- We start with no variables in the model, and as with forward selection, we add the variable that provides the best fit.
- We continue to add variables one-by-one.
- It is possible sometimes that the p-values for variables can become larger as new predictors are added to the model.
- Hence, if at any point the p-value for one of the variables in the model rises above a certain threshold, then we remove that variable from the model.
- We continue to perform these forward and backward steps until all variables in the model have a sufficiently low p-value, and all variables outside the model would have a large p-value if added to the model.

M_0
 $M_1 : \{x_3, f(x_1), f(x_2), f(x_3) \dots f(x_p)\}$

$M_2 : x_3, \min \text{RSS}$

• Stepwise Selection

– Hybrid Approach

- Variables are added to the model sequentially, in analogy to forward selection.
- However, after adding each new variable, the method may also remove any variables that no longer provide an improvement in the model fit.
- Such an approach attempts to more closely mimic best subset selection while retaining the computational advantages of forward and backward stepwise selection.

x_1
 x_2
⋮
 x_p

Linear Model Selection and Regularization

- **Shrinkage**. This approach involves fitting a model involving all p predictors. However, the estimated coefficients are shrunken towards zero relative to the least squares estimates. This shrinkage (also known as regularization) has the effect of reducing variance. Depending on what type of shrinkage is performed, some of the coefficients may be estimated to be exactly zero. Hence, shrinkage methods can also perform variable selection.

Find β_s by $\min RSS$

$$y = \beta_0 + \underbrace{\beta_1}_{\text{A}} x_1 + \underbrace{\beta_2}_{\text{A}} x_2 + \dots + \underbrace{\beta_p}_{\text{A}} x_p$$

Shrinkage

- By retaining a subset of the predictors and discarding the rest, subset selection produces a model that is interpretable and has possibly lower prediction error than the full model.
- However, because it is a discrete process variables are either retained or discarded—it often exhibits high variance, and so doesn't reduce the prediction error of the full model.
- Shrinkage methods are more continuous, and don't suffer as much from high variability.

Shrinkage Methods

1. Ridge Regression

- Ridge regression shrinks the regression coefficients by imposing a penalty on their size. The ridge coefficients minimize a penalized residual sum of squares,

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}.$$

- Here $\lambda \geq 0$ is a complexity parameter that controls the amount of shrinkage. It's directly proportional to shrinkage .

Ridge Regression

$$\underset{\text{Ridge}}{\text{RSS}(\beta)} = \left\| X\beta - Y \right\|_2^2 + \lambda \left\| \beta \right\|_2^2$$

$$\text{RSS}(\beta) = \left\| X\beta - Y \right\|_2^2$$

↑

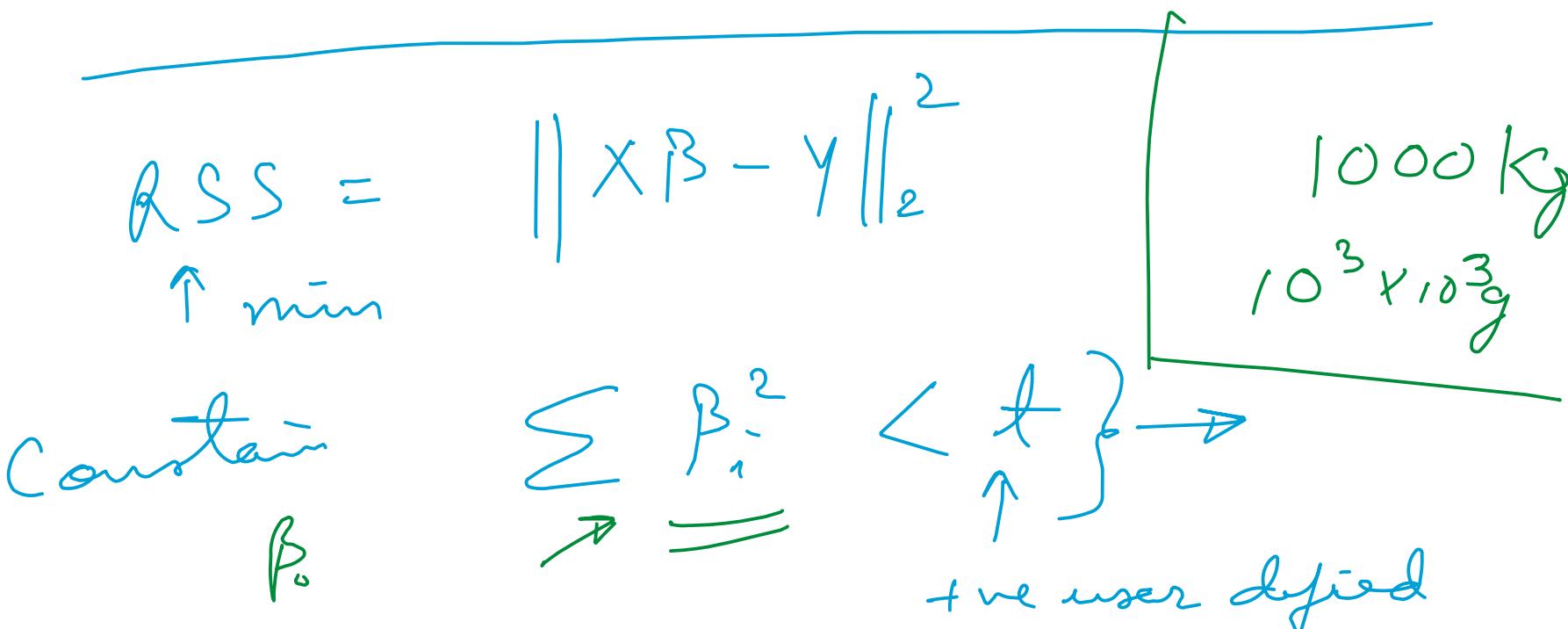
min

constraint $\left\| \beta \right\|_2^2 < t$

↑
+ve

Ridge Regression

$$RSS(\beta) = \text{Ridge} \quad \|X\beta - Y\|_2^2 + \lambda \|\beta\|_2^2$$



Ridge Regression

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 - \dots + \hat{\beta}_p x_p$$



$$\sum \hat{\beta}_i^2 < t$$

$$RSS = \sum_{i=1}^n (\hat{y} - y)^2 = \sum \left(y - \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots \right)^2$$

Ridge Regression

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 - \dots + \hat{\beta}_p x_p$$



$$\sum \hat{\beta}_i^2 < t$$

$$RSS = \sum_{i=1}^n (\hat{y} - y)^2 = \sum \left(y - \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots \right)^2$$

Ridge Regression

$$\min_{\text{Ridge}} \text{RSS}(\beta) = \|X\beta - y\|_2^2 + \lambda \|\beta\|_2^2$$

Tends to 0 as $\lambda \uparrow$

$$\equiv \min \text{RSS}(\beta) = \|X\beta - y\|_2^2$$

Constraint

$$\|\beta\|_2^2 < t$$

- Has no β_0
- x are normalized

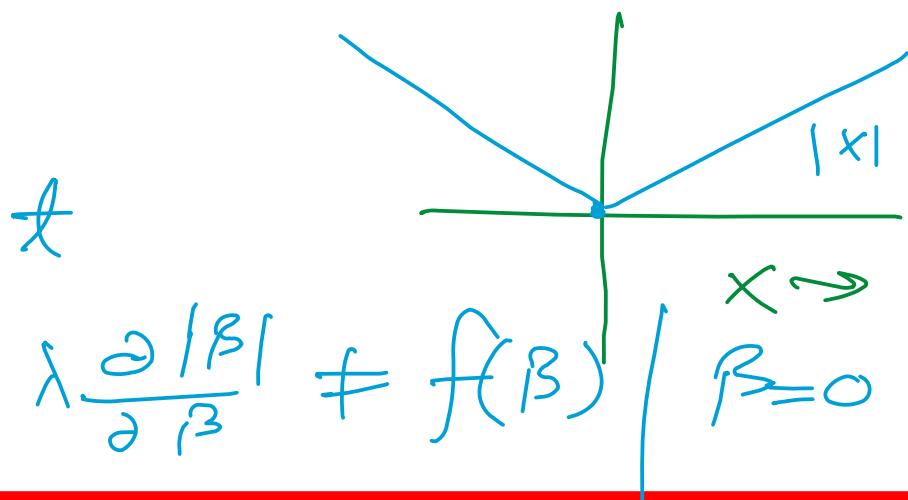
Lasso Regression

$$RSS(\beta) = \text{Ridge} \quad \|X\beta - Y\|_2^2 + \lambda \|\beta\|$$

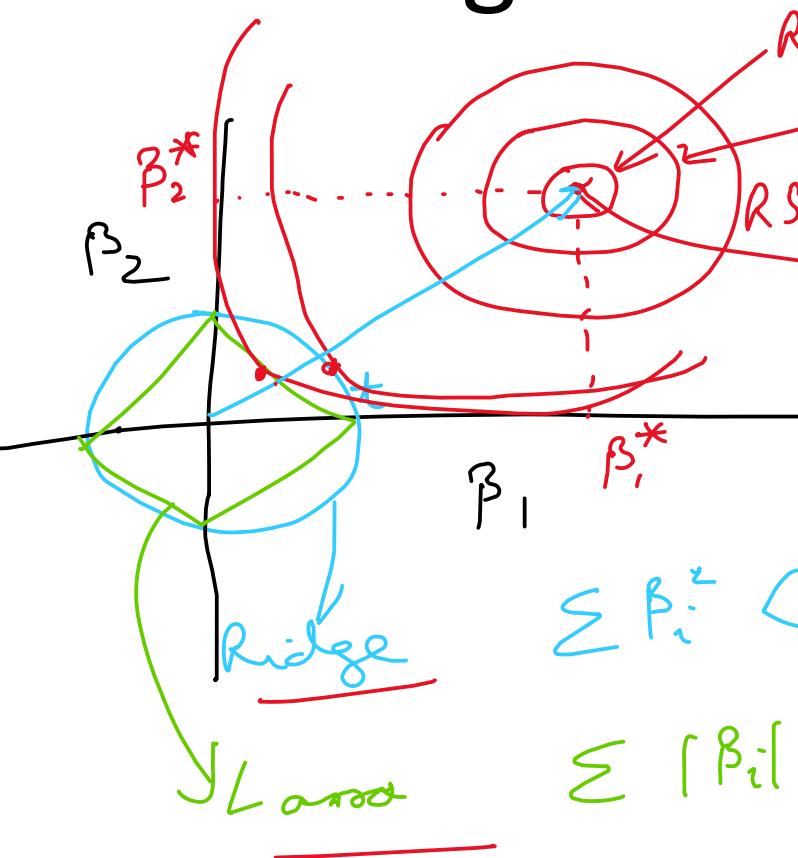
$$\min RSS(\beta) = \|X\beta - Y\|_2^2$$

constraint $\sum |\beta| < t$

$$\lambda \frac{\partial |\beta|}{\partial \beta} = f(\beta)$$
$$\lambda \frac{\partial \beta^2}{\partial \beta} = 2\lambda\beta$$



Ridge and Lasso Regression

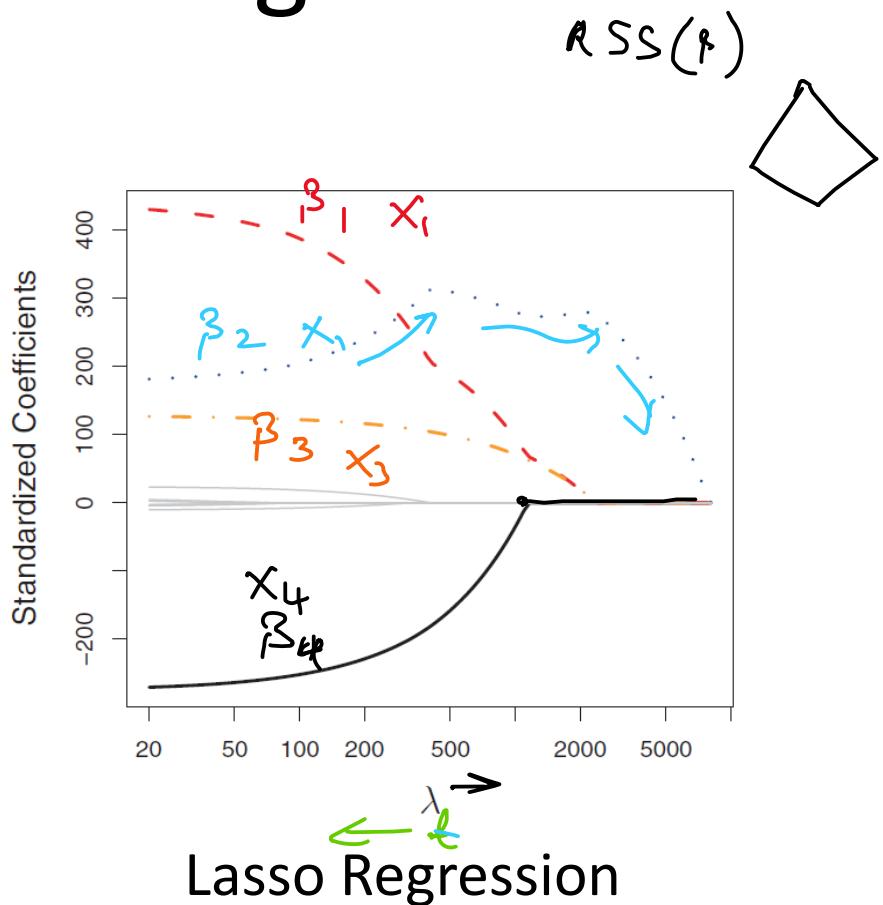
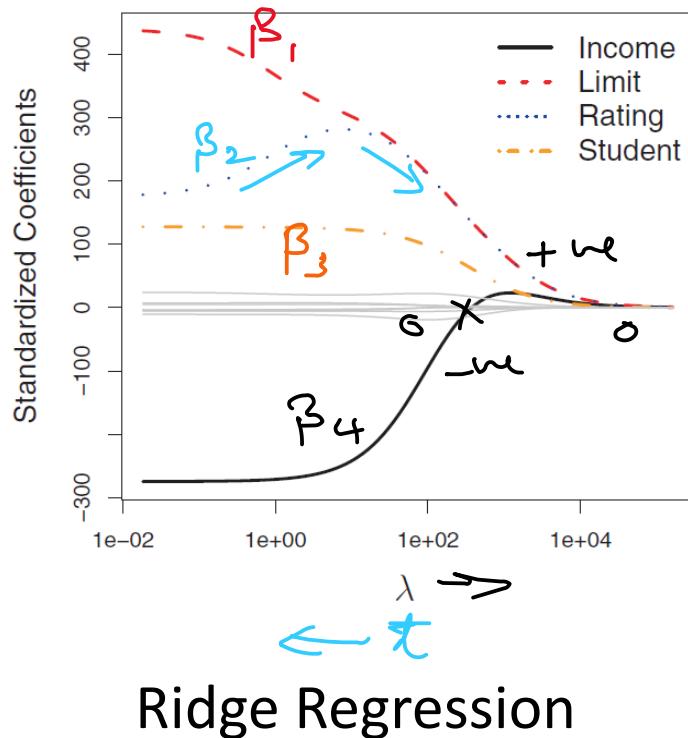


$$\sum \beta_i^2 < t^2$$

$$\sum |\beta_i| < t$$

} minimize RSS
 with the constraint
 $t^2 \geq (\beta_1^*)^2 + (\beta_2^*)^2$
 Ridge is same as simple LR
 t ?

Ridge and Lasso Regression



Ridge and Lasso Regression



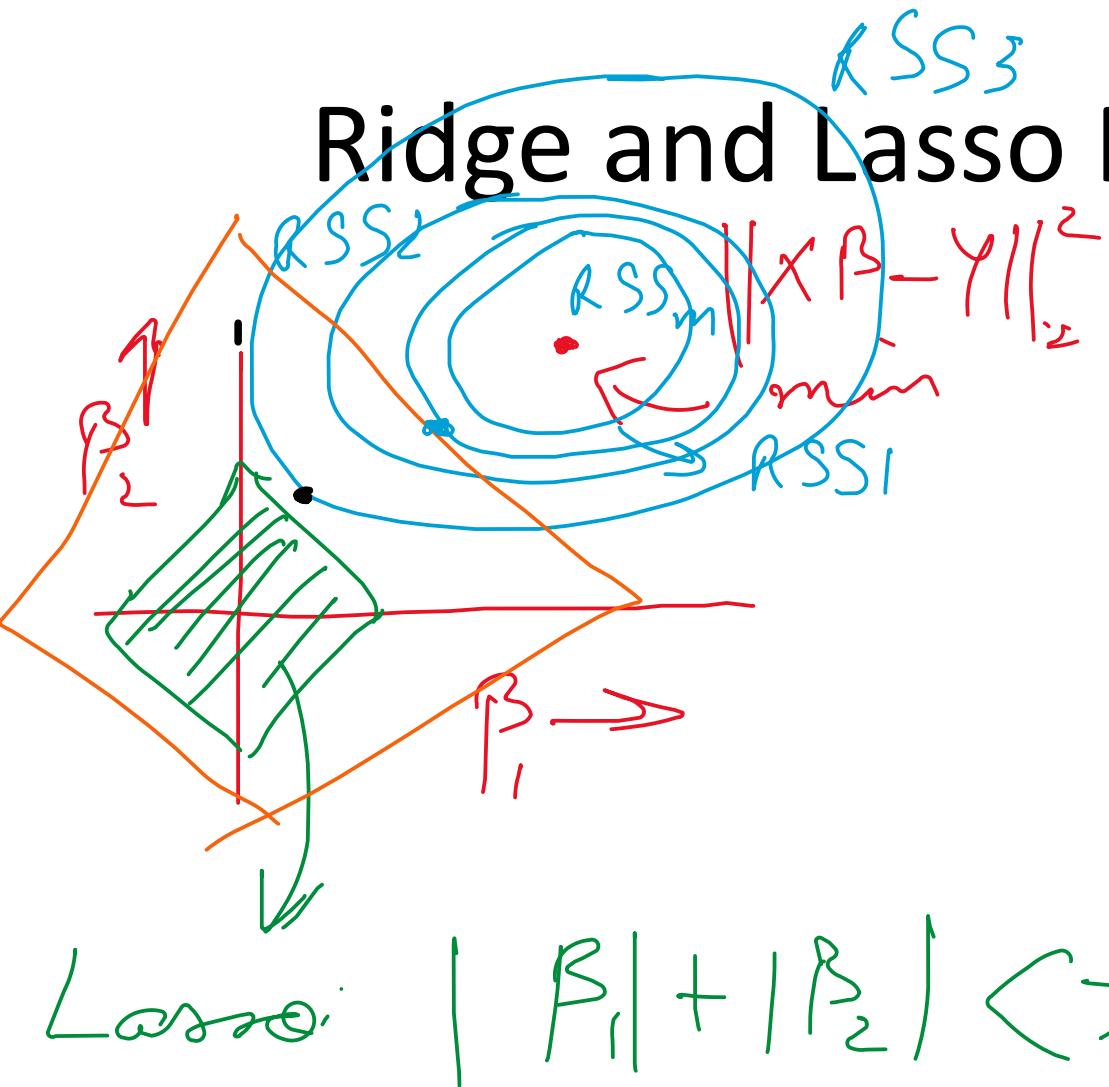
$$RSS_{min} < RSS_1 \\ < RSS_2 < RSS_3$$

Ridge $\sum \beta_i^2 < \alpha$

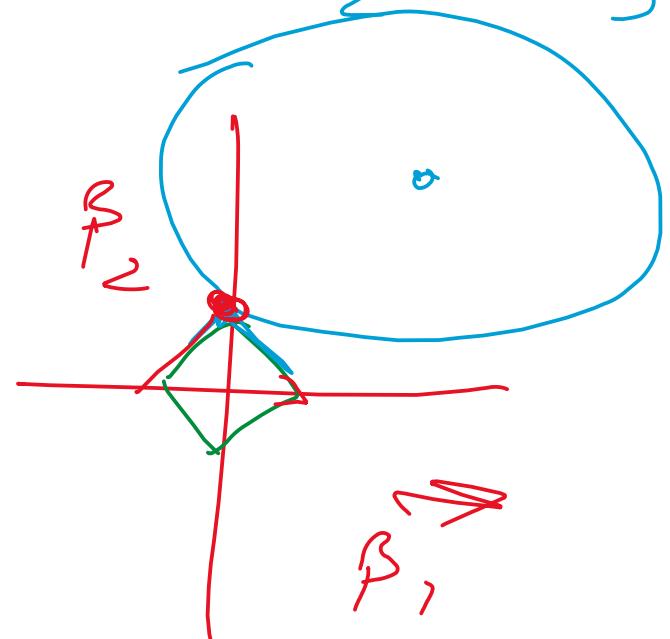
$$\beta_1^2 + \beta_2^2 < \alpha$$

Circle

Ridge and Lasso Regression



$$\begin{aligned} RSS_m &< RSS_1 \\ &< RSS_2 < RSS_3 \end{aligned}$$



Ridge and Lasso Regression

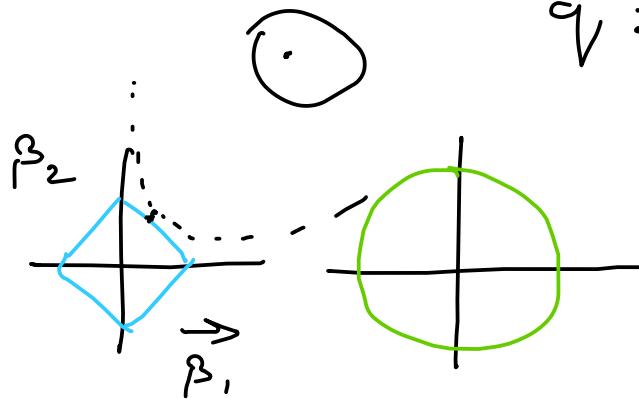
$$RSS(\beta) + \lambda \sum |\beta_i|^q$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

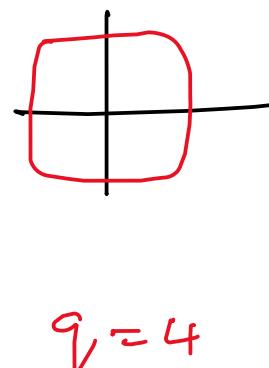
If $q = 0 \rightarrow OLR$

$q = 1$ Lasso

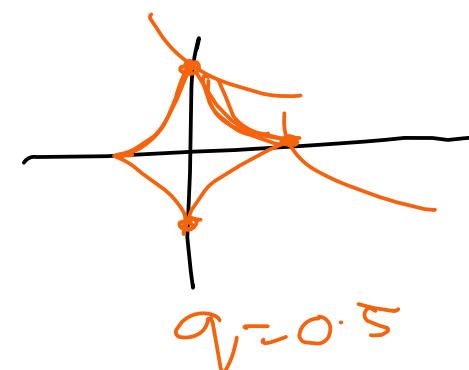
$q = 2$ Ridge



$q = 1$
Lasso



$q = 2$



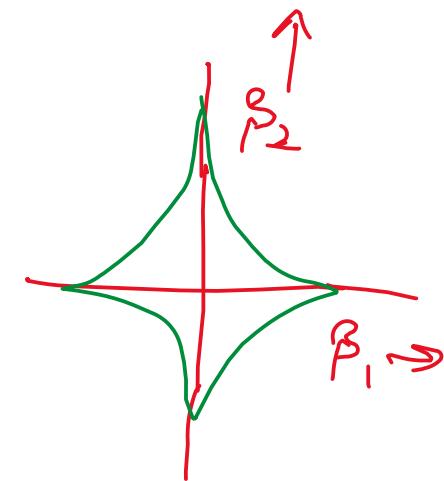
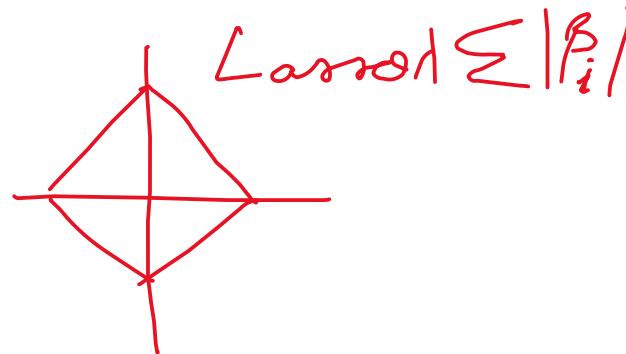
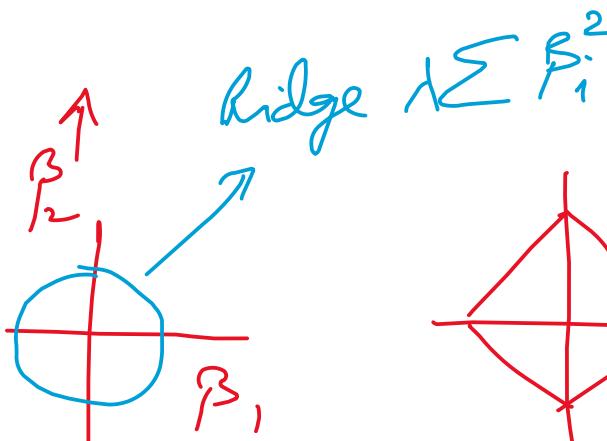
$q = 4$

$q = 0.5$

Ridge and Lasso Regression

General Lasso

$$\lambda \sum |\beta_i|^q$$



$$q=4$$

$$q=2$$

$$q=1$$

$$q = 0.5$$

$$\sum \beta_i^4 < \infty$$

$$|\beta_1|^{0.5} + |\beta_2|^{0.5} < \infty$$

Shrinkage Methods

- An equivalent way to write the ridge problem is

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2,$$

$$\text{subject to } \sum_{j=1}^p \beta_j^2 \leq t,$$

Ridge Regression

- Ridge regression is very similar to least squares, except that the coefficients ridge regression are estimated by minimizing a slightly different quantity.

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 = \text{RSS} + \lambda \sum_{j=1}^p \beta_j^2,$$

- where $\lambda \geq 0$ is a tuning parameter. The second term, is called a shrinkage penalty, it is small when β_1, \dots, β_p are close to zero

Shrinkage Methods

2.The lasso

Lasso is a shrinkage method like ridge, with subtle but important differences. The lasso estimate is defined by

$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2$$

subject to $\sum_{j=1}^p |\beta_j| \leq t.$

Shrinkage Methods

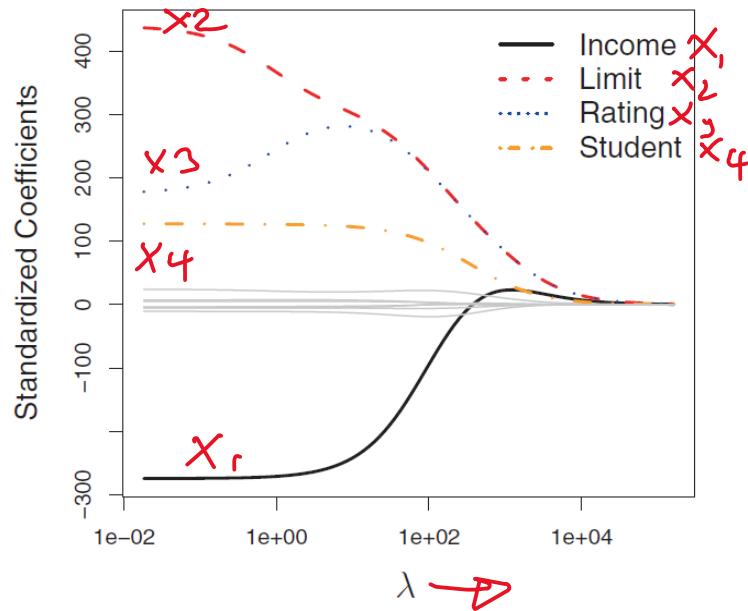
- We can also write the lasso problem in the equivalent Lagrangian form

$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2} \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij}\beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}.$$

$$\frac{\partial (\text{RSS})_{\text{Lasso}}}{\partial \beta_i} = \dots + \overset{\circ}{\lambda} \quad \text{Lasso}$$

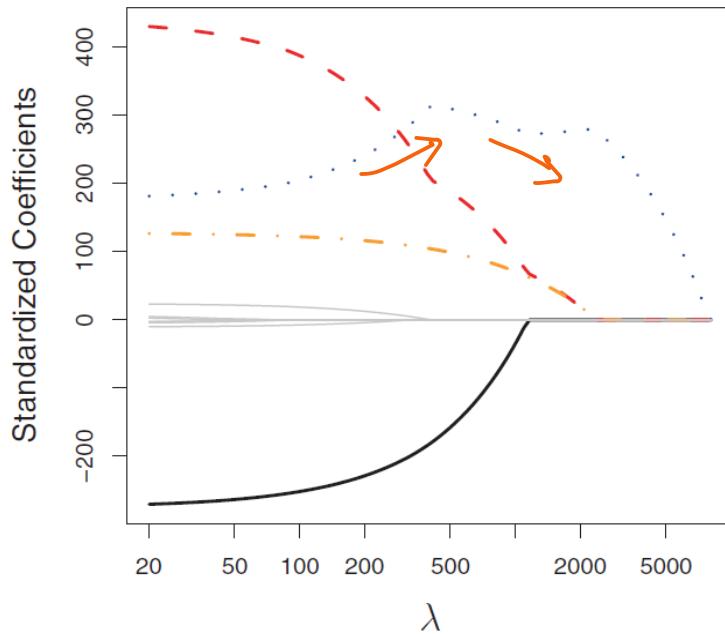
$$\frac{\partial (\text{RSS})_{\text{Ridge}}}{\partial \beta_i} = \dots + \overset{\circ}{2\lambda\beta_i}$$

Ridge and Lasso Regression



Ridge Regression

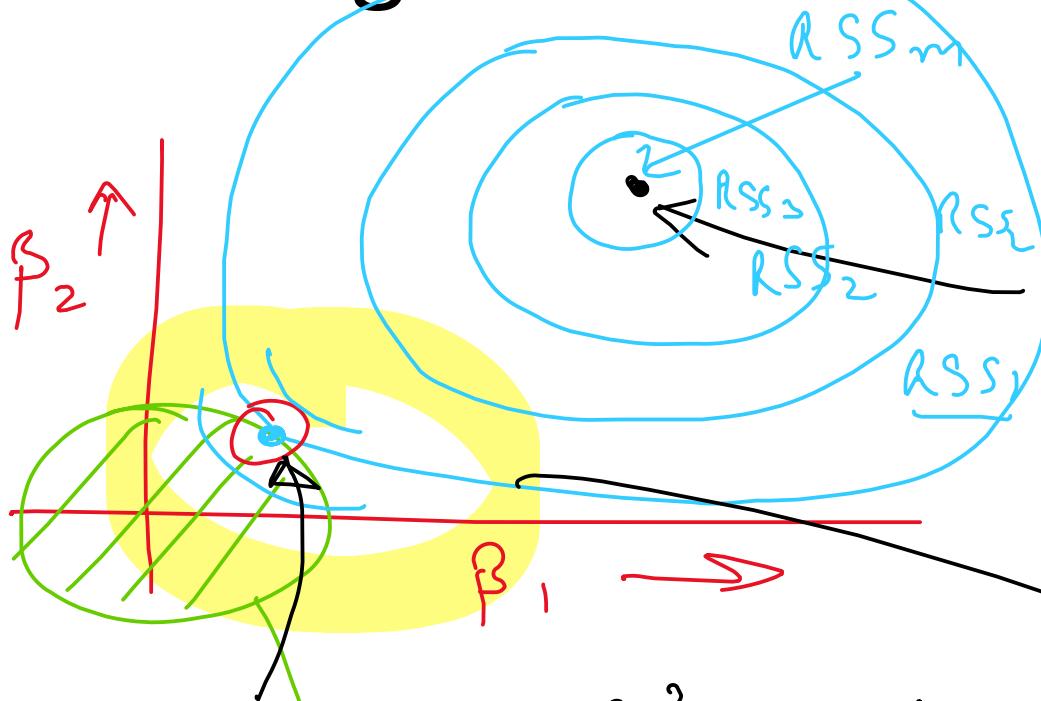
$$\lambda \sum \beta_i^2$$



Lasso Regression

$$\lambda \sum |\beta_i|$$

Ridge and Lasso Regression



Ridge :

$$\sum \beta_i^2 < t$$

$$\beta_1^2 + \beta_2^2 < t$$

$$RSS = \| X\beta - Y \|_2^2$$

\min

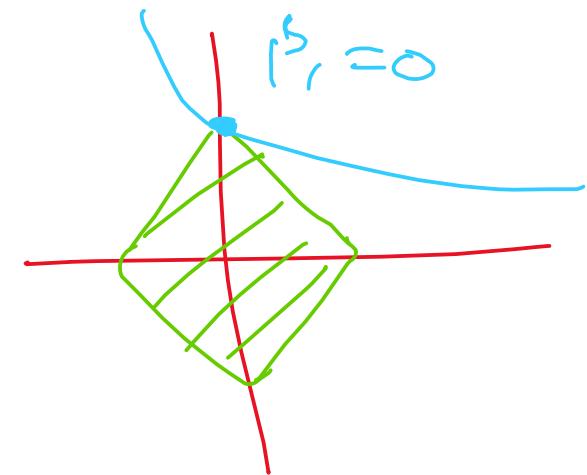
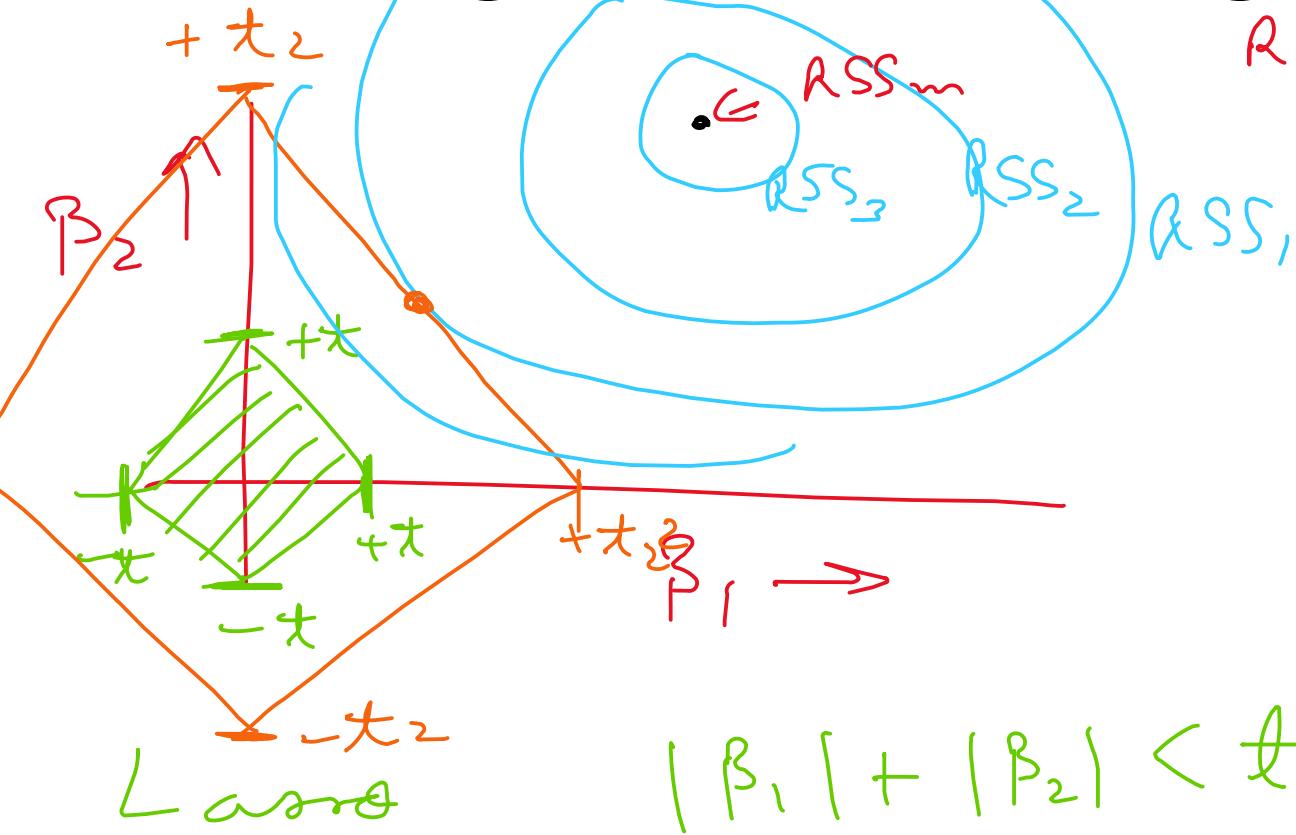
$\min RSS$

$$RSS_m < RSS_3 < RSS_2 \\ < RSS_1$$

Tangent
is the sol.
for Ridge Regress.

Ridge and Lasso Regression

$$RSS = \|X\beta - Y\|_2^2$$



Summary

- Ridge regression does a proportional shrinkage.
- Lasso translates each coefficient by a constant factor , truncating at zero. This is called “soft thresholding.”.
- Best-subset selection drops all variables with coefficients smaller than the Mth largest; this is a form of “hard-thresholding.”

Dimensionality reduction

- Feature selection: Feature selection approaches try to find a subset of the original variables (also called features or attributes)
 - Filter strategy (e.g. information gain)
 - Wrapper strategy (e.g. search guided by accuracy)
 - Embedded strategy (features are selected to add or be removed while building the model based on the prediction errors)
- Feature projection: Feature projection transforms the data in the high-dimensional space to a space of fewer dimensions. The data transformation may be linear, as in principal component analysis (PCA), but many nonlinear dimensionality reduction techniques also exist. For multidimensional data, tensor representation can be used in dimensionality reduction through multilinear subspace learning.

Feature projection

- Principal component analysis (PCA)
- Non-negative matrix factorization (NMF)
- Kernel PCA
- Graph-based kernel PCA
- Linear discriminant analysis (LDA)
- Generalized discriminant analysis (GDA)
- Autoencoder

Principal Component Analysis

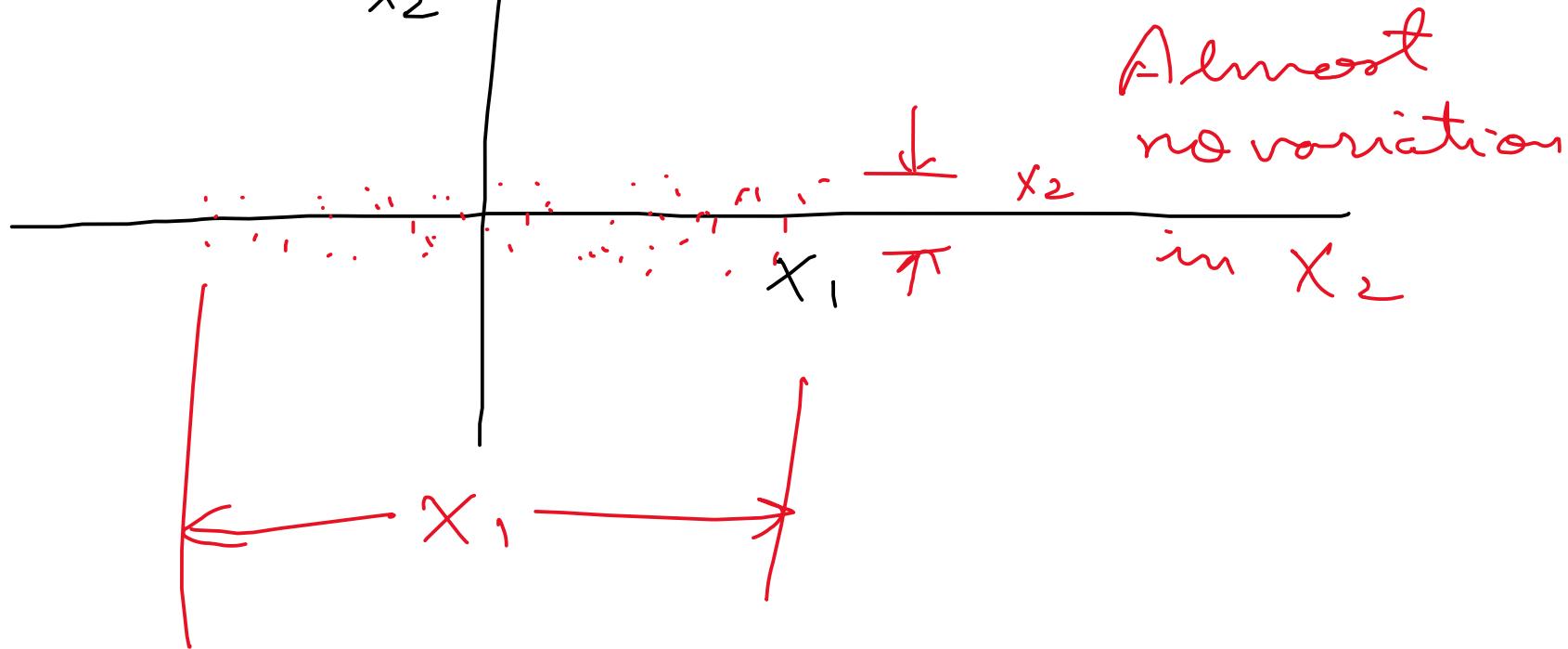
$$y = f(x_1, x_2) \simeq f(x_1)$$



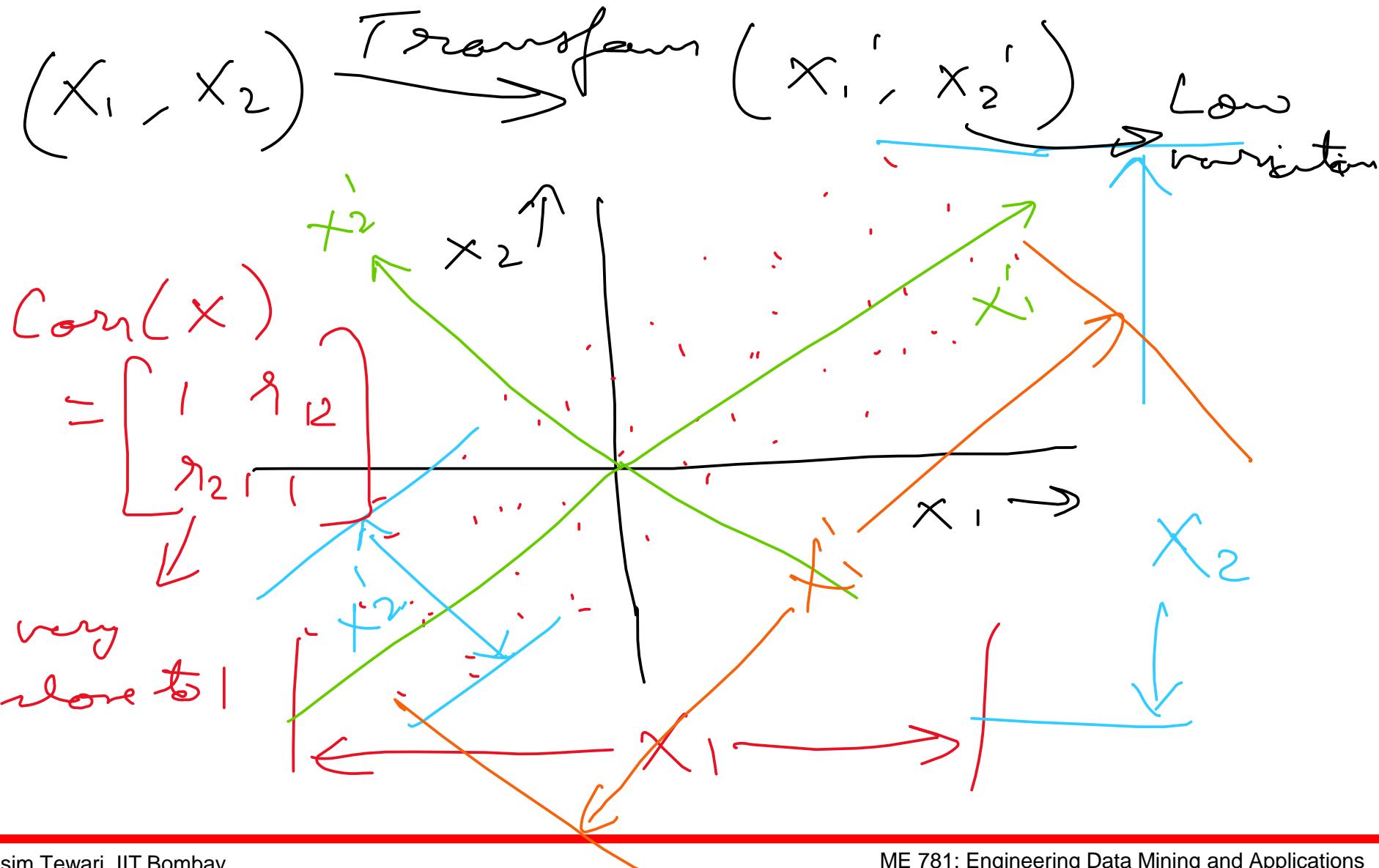
x_2

$$f(x_1)$$

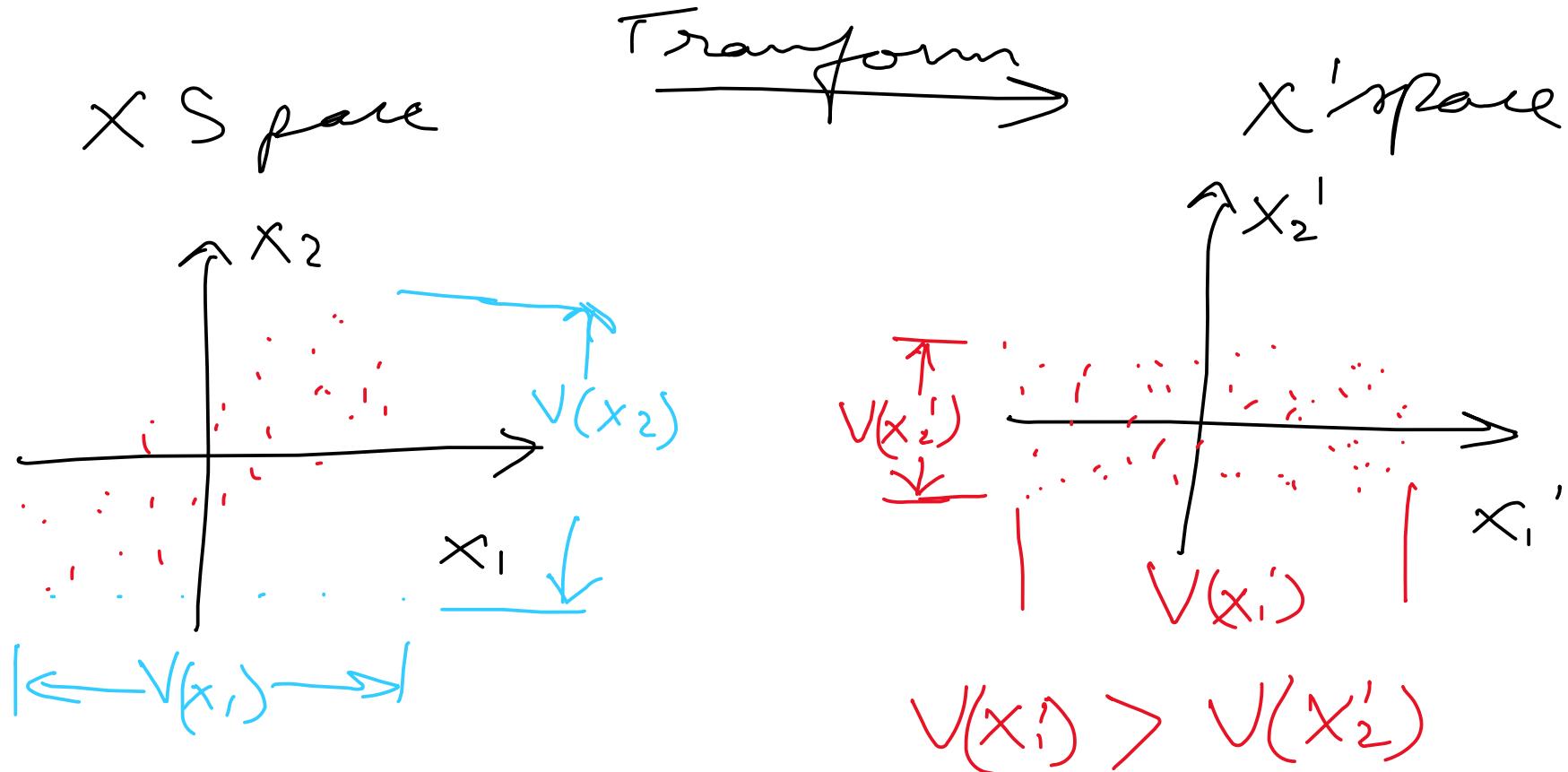
No variation



Principal Component Analysis



Principal Component Analysis



Principal Component Analysis

In 2D space with n data points

X space \rightarrow X' space

$$X = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{n1} & x_{n2} \end{bmatrix} \Rightarrow X' = \begin{bmatrix} x'_1 & x'_{12} \\ x'_{21} & x'_{22} \\ \vdots & \vdots \\ x'_{n1} & x'_{n2} \end{bmatrix}$$

Variance

$$V(x'_1) > V(x'_2)$$

↑
highest possible variance

Principal Component Analysis

In P -dimensional Space

PCA
 x space $\rightarrow x'$ space

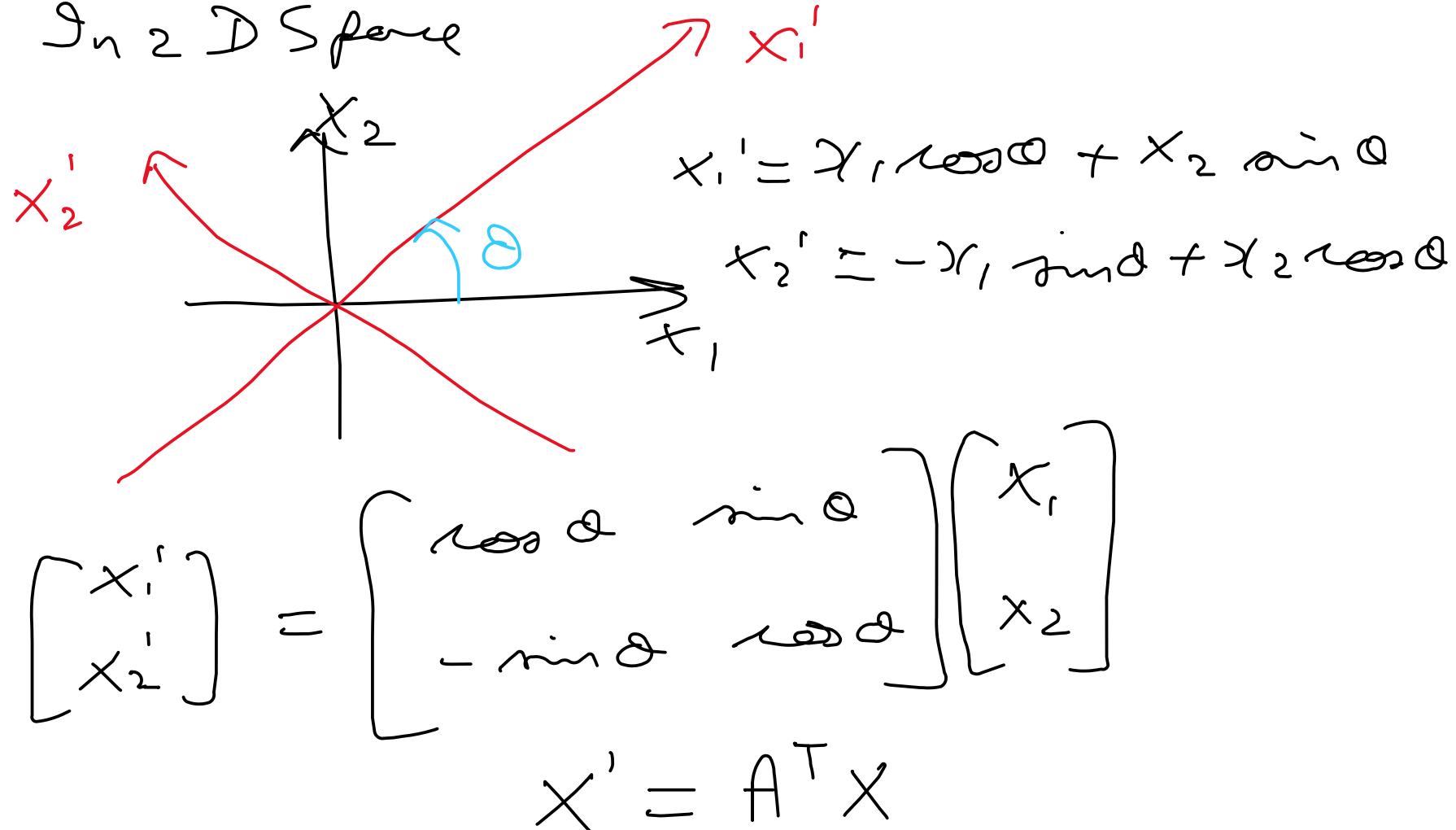
$$X_{n \times p} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}$$

$$X'_{n \times p} = \begin{bmatrix} x'_1 & x'_{12} & \dots & x'_{1p} \\ x'_{21} & x'_{22} & \dots & x'_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x'_{n1} & x'_{n2} & \dots & x'_{np} \end{bmatrix}$$

$$V(x'_1) \geq V(x'_2) \geq \dots \geq V(x'_p)$$

Principal Component Analysis

In 2 D Space



Principal Component Analysis

In 3D

$$x' = A^T x$$

$$\begin{bmatrix} x_1' \\ x_2' \\ x_3' \end{bmatrix} = \begin{bmatrix} \cos \theta & \sin \theta & 0 \\ -\sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$R_3(\theta) = A^T \quad \text{For rotation about } x_3$$

Principal Component Analysis

$$R_2(\theta) = \begin{bmatrix} \cos \theta & 0 & \sin \theta \\ 0 & 1 & 0 \\ -\sin \theta & 0 & \cos \theta \end{bmatrix}$$

$$R_1(\phi) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos \phi & \sin \phi \\ 0 & -\sin \phi & \cos \phi \end{bmatrix}$$

For a general rotation

$$A^T = R = R_1(\alpha) R_2(\beta) R_3(\gamma)$$

Principal Component Analysis

In general 3D rotation you can always find an axis that is not changed.

$$x' = R x$$

↑
new
axis

↑
old
axis

Invariant axis

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\begin{bmatrix} x \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} = R \begin{bmatrix} x \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \Rightarrow (R - I)x = 0$$

Invariant vector

Principal Component Analysis

$$X' = A^T X$$

$$P \times 1 \quad P \times P \quad P \times 1$$

In 2D

$$A^T = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}$$

$$\Rightarrow A = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} = [a_1 \ a_2]$$

$$a_1 = \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} ; a_2 = \begin{bmatrix} -\sin \theta \\ \cos \theta \end{bmatrix}$$

Principal Component Analysis

$$A = \begin{bmatrix} a_1 & a_2 \\ 2 \times 2 & 2 \times 1 \end{bmatrix} \quad \left| \quad \begin{array}{l} a_1 = \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} \\ a_2 = \begin{bmatrix} -\sin \theta \\ \cos \theta \end{bmatrix} \end{array} \right.$$

$$a_1^T a_1 = [\cos \theta \ \sin \theta] \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} = \cos^2 \theta + \sin^2 \theta = 1$$

$$a_2^T a_2 = 1$$

$$a_1^T a_2 = [\cos \theta \ \sin \theta] \begin{bmatrix} -\sin \theta \\ \cos \theta \end{bmatrix} = 0$$

Principal Component Analysis

$$A^T A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I = A A^T = A A^{-1}$$

$$\Rightarrow A^T = A^{-1}$$

$\therefore A$ is an orthogonal Transformation

Principal Component Analysis

In P-dimension

$$A = [a_1 \ a_2 \ \dots \ a_p]$$

$$a_i^T a_j = S_{ij} = \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j \end{cases}$$

In PCA we want this transformation

to be such that variances of x' are as follows:

$$V(x'_1) \geq V(x'_2) \geq V(x'_3) \dots \geq V(x'_P)$$

Principal Component Analysis

p-dimensional Space

$$X_{n \times p} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ \vdots & & & \\ x_{k1} & x_{k2} & \dots & x_{kp} \\ \vdots & & & \\ x_{n1} & \dots & \dots & x_{np} \end{bmatrix}_{n \times p}$$

In set form $X = \{x_1, x_2, \dots, x_n\}$

where $x_k = (x_{k1}, x_{k2}, \dots, x_{kp})$

Principal Component Analysis

∴ for the "kth" Set

$$x_k = \begin{bmatrix} x_{k1} \\ x_{k2} \\ \vdots \\ x_{kp} \end{bmatrix}^T$$

Similarly

$$\begin{bmatrix} x'_{k1} \\ x'_{k2} \\ \vdots \\ x'_{kp} \end{bmatrix}^T = \begin{bmatrix} x'_{k1} & x'_{k2} & \dots & x'_{kp} \end{bmatrix} = x'_k$$

$$\begin{bmatrix} x'_{k1} \\ x'_{k2} \\ \vdots \\ \vdots \\ x'_{kp} \end{bmatrix} = A^T \begin{bmatrix} x_{k1} \\ x_{k2} \\ \vdots \\ \vdots \\ x_{kp} \end{bmatrix}$$

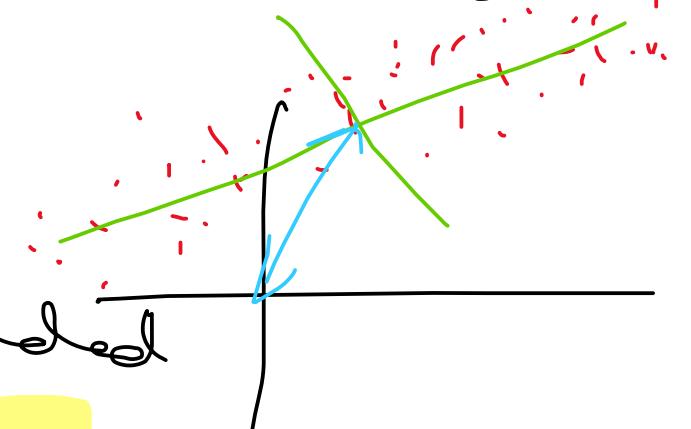
Principal Component Analysis

$$\begin{bmatrix} x'_{k_1} \\ x'_{k_2} \\ \vdots \\ x'_{k_p} \end{bmatrix}^T = A^T \begin{bmatrix} x_{k_1} \\ x_{k_2} \\ \vdots \\ x_{k_p} \end{bmatrix}^T$$

$$\begin{bmatrix} x'_{k_1} \\ x'_{k_2} \\ \vdots \\ x'_{k_p} \end{bmatrix}^T = \begin{bmatrix} x_{k_1} \\ x_{k_2} \\ \vdots \\ x_{k_p} \end{bmatrix}^T A$$
$$x'_{\underset{\uparrow}{k}} = \underset{\uparrow}{x_k} A$$

Principal Component Analysis

$$x'_k = x_k A$$



Translation is also needed

$$x'_k = (x_k - \bar{x}) A$$

$1 \times P$ $1 \times P$ $P \times P$

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k = \frac{1}{n} \left[\sum x_{k1}, \sum x_{k2}, \dots, \sum x_{kP} \right]$$

$$x_k = x'_k A^T + \bar{x}$$

Principal Component Analysis

$$x_k = x'_k A^T + \bar{x}$$

Variance vector of $x' = \begin{bmatrix} \text{var}(x'_1) & \text{var}(x'_2) \\ & \dots \text{var}(x'_p) \end{bmatrix}$

Covariance of $x' = \Sigma_{x'} = \frac{1}{n-1} \sum_{k=1}^n (x'_k)^T (x'_k)$

$\uparrow \quad \uparrow \quad \uparrow$
 $P \times P \quad P \times 1 \quad 1 \times P$

$$\Sigma_{x'} = \frac{1}{n-1} \sum_{k=1}^n (x'_k)^T (x'_k)$$

Principal Component Analysis

$$\Sigma_{x'} = \frac{1}{n-1} \sum_{k=1}^n (x'_k)^T (x'_k)$$

$$x'_k = (x_k - \bar{x}) A$$

$$\Sigma_{x'} = \frac{1}{n-1} \sum \left[(x_k - \bar{x}) A \right]^T \left[(x_k - \bar{x}) A \right]$$

$$= \frac{1}{n-1} \sum_{k=1}^n A^T (x_k - \bar{x})^T \cdot (x_k - \bar{x}) A$$

$$= A^T \left[\frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})^T \cdot (x_k - \bar{x}) \right] A$$

$$\Sigma_{x'} = A^T C_A A$$

Covariance matrix of x

Principal Component Analysis

$$\mathcal{V}_{x^1} = A^T \underbrace{C_A}_\text{Covariance matrix of } X A$$

$$C_{ij} = \frac{1}{n-1} \sum_{k=1}^n \left[x_{ki} - \frac{1}{n} \sum_{l=1}^n x_{li} \right] \left[x_{kj} - \frac{1}{n} \sum_{l=1}^n x_{lj} \right]$$

$$i, j \in \{1, 2, \dots, p\}$$

Choose A such that \mathcal{V}_{x^1} is maximized.

However $A^T A = I$

Principal Component Analysis

Maximise λ_x , given $A^T A = I$

Constrained optimization.

Construct a Lagrange function

$$L = A^T C A - \lambda (A^T A - I)$$

Sol is A which maximises L

$$\begin{aligned} \Rightarrow \frac{\partial L}{\partial A} = 0 &\Rightarrow C A - \lambda A = 0 \\ &\Rightarrow (C - \lambda I) A = 0 \end{aligned}$$

Principal Component Analysis

$$(C - \lambda I) A = 0 \quad \text{Solve for } A.$$

This is an Eigenvalue/Eigenvector Problem.

This started by solving DEs:

$$\frac{dy}{dt} = Ay \quad \left. \begin{array}{l} \\ \uparrow \\ n \times n \end{array} \right\} n \text{ linear DEs}$$

Sol is of the type $y(t) = e^{\lambda t} X$

Principal Component Analysis

$$y(t) = e^{\lambda t} x$$

$$\Rightarrow \lambda e^{\lambda t} x = A e^{\lambda t} x$$

$$\Rightarrow \lambda x = A x$$

$$(A - \lambda I)x = 0$$

$n \times n$ $n \times n$ $n \times 1$

$n \times 1$

$$\frac{dy}{dt} = A y \quad \left. \begin{array}{l} \text{n linear} \\ \text{DEs} \end{array} \right\}$$

$$(A - \lambda I) A = 0$$

$P \times P$ $P \times P$ $P \times P$

$P \times P$

$$(A - \lambda I)x = 0$$

\uparrow
 $n \times 1$

} get eigenvalues of λ
as eigenvectors of X

Eq.

$$\begin{aligned} y_1' &= 5y_1 + y_2 \\ y_2' &= 3y_1 + 3y_2 \end{aligned} \quad \left\{ \frac{dy}{dt} = \begin{bmatrix} 5 & 1 \\ 3 & 3 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right.$$

$$\lambda_1 = 6 ; x_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} ; \lambda_2 = 2 ; x_2 = \begin{bmatrix} 1 \\ -3 \end{bmatrix}$$

$$y(t) = \begin{bmatrix} y_1(t) \\ y_2(t) \end{bmatrix} = c_1 e^{\lambda_1 t} x_1 + c_2 e^{\lambda_2 t} x_2$$

$$y(t) = c_1 e^{6t} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + c_2 e^{2t} \begin{bmatrix} 1 \\ -3 \end{bmatrix}$$

↑ get from BC

~~Case 1~~ If we know the eigenvector for
 $Ax = \lambda_1 x$, then we also know
the eigenvector for $A^2x = \lambda_2 x$

$$A^2x = \lambda_2 x \Rightarrow A(Ax) = \lambda_2 x$$

$$\Rightarrow A(\lambda_1 x) = \lambda_2 x$$

$$\Rightarrow \lambda_1^2 x = \lambda_2 x$$

$$\Rightarrow \lambda_2 = \lambda_1^2$$

and they have the same eigenvector

Case 2

$$\text{If } (A + cI)x = \lambda_3 x$$

$$\lambda_3 = (\lambda_1 + c)$$

Similarly $A^n x = \lambda_n x$

then $\lambda_n = (\lambda_1)^n$

$$(C - \lambda I)A = 0$$

$P \times P$

This is a concatenation of eigenvectors

$$(A - \lambda I)x = 0$$

$n \times 1$

}

$$A = (a_1, a_2, \dots, a_p)$$

\uparrow

$P \times P$

$$(C - \lambda I) A = 0$$

$$A = (\alpha_1, \alpha_2 \dots \alpha_p)$$

↑ ↑

eigenvectors \leftrightarrow eigenvalues

\Rightarrow Variance in x' corresponds to the eigenvalue of $\lambda_1, \lambda_2, \dots, \lambda_p$

$$CA = \lambda A \Rightarrow A^T C A = A^T \lambda A$$

$$\lambda = A^T C A \equiv \Sigma_{x_i}$$

$$\lambda_1 > \lambda_2 > \lambda_3 \dots > \lambda_p$$

Principal Component Analysis

Thus, if we want to capture 95% variance, then we reduce the # of dimensions to q such that

$$\frac{\sum_{i=1}^q \lambda_i}{\sum_{i=1}^p \lambda_i} \geq \frac{95}{100}$$

Principal Component Analysis

Eg. $P = 2, n = 4$

$$X = \{(1,1), (2,1), (2,2), (3,2)\}$$

$$\Rightarrow \bar{x} = \left(2, \frac{3}{2}\right); C = \frac{1}{3} \begin{bmatrix} 2 & 1 \\ 1 & 1 \end{bmatrix}$$

The eigenvalues are

$$\lambda_1 = 0.8727, a_1 = \begin{bmatrix} -0.85065 \\ -0.52571 \end{bmatrix}$$

$$\lambda_2 = 0.1273, a_2 = \begin{bmatrix} 0.52571 \\ -0.85065 \end{bmatrix}$$

$$A = [a_1, a_2] = \begin{bmatrix} -0.85065 & 0.52571 \\ -0.52571 & -0.85065 \end{bmatrix}$$

Principal Component Analysis

$$x' = A^T(x - \bar{x}) \Rightarrow x_1' = \begin{bmatrix} -0.85085 & -0.52571 \\ 0.52571 & -0.85085 \end{bmatrix} \begin{bmatrix} 1 & -2 \\ 1 & -1.5 \end{bmatrix}$$

$$\Rightarrow \underline{x_1}' = \begin{bmatrix} 1.1135 \\ -0.10039 \end{bmatrix}$$

$$\underline{x_2}' = \begin{bmatrix} 0.2628 \\ 0.42533 \end{bmatrix}; \quad \underline{x_3}' = \begin{bmatrix} -0.2628 \\ -0.42533 \end{bmatrix}$$

$$\underline{x_4}' = \begin{bmatrix} -1.1135 \\ 0.10039 \end{bmatrix}$$

$x \xrightarrow[A^T]{\text{Transformed}} x'$

Principal Component Analysis

$$X \xrightarrow{A^T} X'; \quad A^T = \begin{bmatrix} -0.85065 & -0.52571 \\ 0.52571 & -0.85065 \end{bmatrix}$$

X

$\{(1,1), (2,1), (2,2), (3,2)\}$

X'

$\{(1.1135, -0.10039), (0.26286, 0.42533), (-0.26286, -0.42533), (-1.1135, -0.10039)\}$

mean $X = (2, 1.5)$

Covariance of X

$$= \begin{bmatrix} \frac{2}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} \end{bmatrix}$$

Mean $X' = (0, 0)$

Cov $X' = \begin{bmatrix} 0.87266 & 0 \\ 0 & 0.12732 \end{bmatrix}$

Principal Component Analysis

$$X \xrightarrow{A^T} X'$$

$$\text{Cov } X' = \begin{bmatrix} \lambda_1 & & & 0 \\ 0 & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_p \end{bmatrix}$$

$$X \xleftarrow{A} X'$$
$$X_1, X_2 \xrightarrow{A^T} X'_1, X'_2$$

Since X'_1 contains 87% of variance, we can drop X'_2

Principal Component Analysis

Let us drop x_2'

Thus we do not take $A = [a, a_2]$

but instead $A' = [a, 0]$

$$X \xrightarrow{(A')^T} X'' \xrightarrow{A'} \sim X$$
$$(A')^T = \begin{bmatrix} -0.85065 & -0.52571 \\ 0 & 0 \end{bmatrix}$$

$$X = \begin{Bmatrix} (1,1) \\ (2,1) \\ (2,2) \\ (3,2) \end{Bmatrix} \quad X'' = \begin{Bmatrix} (1.11351, 0) \\ (0.26286, 0) \\ (-0.26286, 0) \\ (-1.11351, 0) \end{Bmatrix} \xrightarrow{A'} \begin{Bmatrix} (1.0528, 0.91412) \\ (1.7764, 1.36181) \\ (2.2236, 1.63819) \\ (2.9472, 2.08538) \end{Bmatrix}$$

Principal Component Analysis

$$X = \{(1, 1), (2, 1), (2, 2), (3, 2)\}$$

$$\bar{x} = \frac{1}{2} \cdot (4, 3)$$

$$C = \frac{1}{3} \cdot \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix}$$

$$\lambda_1 = 0.8727$$

$$\lambda_2 = 0.1273$$

$$v_1 = \begin{pmatrix} -0.85065 \\ -0.52573 \end{pmatrix}$$

$$v_2 = \begin{pmatrix} 0.52573 \\ -0.85065 \end{pmatrix}$$

$$E = v_1 = \begin{pmatrix} -0.85065 \\ -0.52573 \end{pmatrix}$$

The projected data

$$Y = \{1.1135, 0.2629, -0.2629, -1.1135\}$$

Inverse PCA yields

$$X' = \{ (1.0528, 0.91459), (1.7764, 1.3618), (2.2236, 1.6382), (2.9472, 2.0854) \} \neq X$$

