

Midterm Exam

Start time : 12:50 PM CST 3/10/2023

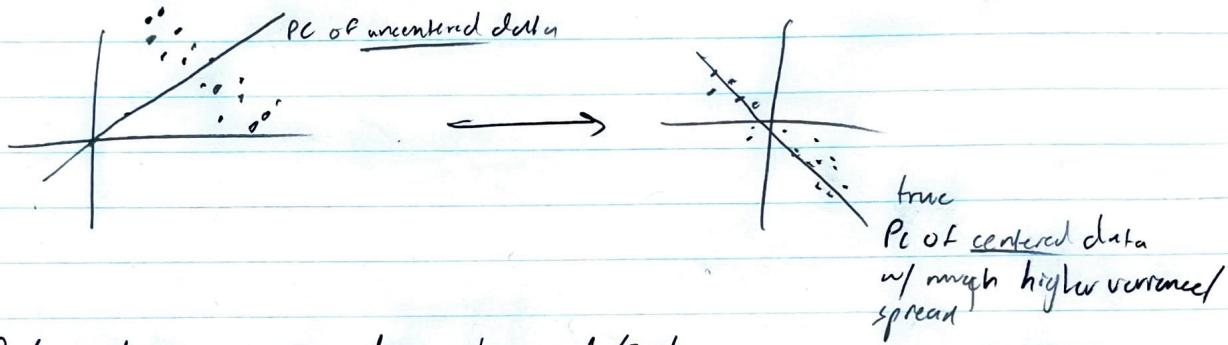
End time : 3:50 PM CST

1a. Sparse vector \mathbb{R}^5 :

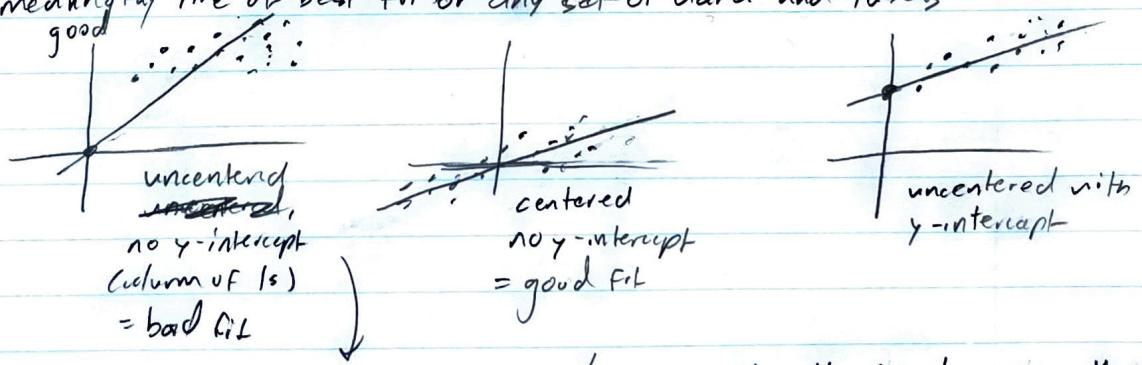
$$\mathbf{x} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 3 \\ 0 \end{bmatrix}$$

most elements 0

b. Data centering is important for accurately identifying the real principal components of a data set since the 1D subspace must pass through the origin:



c. Data centering or appending column of 1s to a data matrix to account for a y-intercept is important to produce a meaningful line of best fit for any set of data and labels



without centering or a column of 1s/y-intercept, the trend within the data set cannot be captured as well.

d. Ridge or lasso regression are forms of regularization to account for singular or ill-conditioned data sets where a unique set of coefficients (w) cannot be found or is very "unstable" (in the case of very small but non-zero singular values). The penalty terms add a constraint to the optimization/regression to account for this and produce a stable/unique set of coefficients w .

Continued on next page

(1). (d). Ridge regression has an analytical equation that can be solved

$$X^T y = (X^T X + \lambda I) w_{\text{ridge}}$$

since it uses the 2-norm, while lasso reg. must be solved with a more complex optimization algorithm, which can be more complex/less intuitive.

Lasso regression produces a more ~~sparse~~^{sparse} set of coefficients w , effectively helping with variable selection (selecting the most impactful variables and "ignoring" [0-weight] the others) while ridge does not which may be good for simplifying the number of vars. in the problem

c. "Linear" in lin. reg. refers to the fact that the terms in the model depend linearly on the weighty coefficients w — i.e. in the line of best fit case in 2D:

$$\hat{y} = w_0 + w_1 x$$

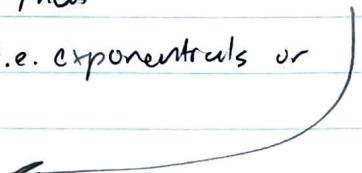
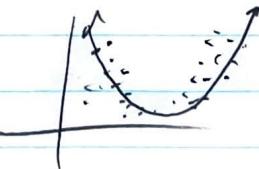
The actual model ~~is~~ can be nonlinear w's are multiplied to x's
however (i.e. the dependence of y prediction (~~is~~ linear combination)
on the data, x) — for quadratics:

$$\hat{y} = w_0 + w_1 x + w_2 x^2$$

$\downarrow x^2 \text{ term}$

This nonlinearity can be useful for data sets that follow a clearly non-linear relationship (i.e. exponentials or ~~quadratics~~ quadratics):

where a ~~line~~ line of best fit is not ideal



f. Overfitting refers to the case in linear regression where there are more parameters w ; than data points (x_i, y_i) , which causes the model to be fit to ~~the data~~ worse in the data. This occurs when $n < p$ for the data matrix — i.e. the matrix has more columns than rows and is short & wide, causing some columns of X to be linearly independent. \rightarrow vectors

g. Next page

(1). g. If a function is convex for its entire domain, then it must have a single / unique global minimum. ~~The gradient~~ The gradient in this case always points "downhill" towards the global minimum from any point in the domain of the function.

h. Gradient descent uses the entire set of objective functions each iteration/step (= epoch); the updated step is $w^{t+1} = w^t - \mu^t \sum_{f=1}^T \nabla L_f(w)$ using all T objective functions $L_f(w)$ in the step.

Stochastic gradient descent uses a ~~set~~ ^{random} subset of the objective functions each iteration/step until all objectives have been used (1 epoch); the step now is $w^{t+1} = w^t - \mu^t \sum_{f \in T'} \nabla L_f(w)$ where T' is a set of randomly selected gradient components that is chosen for the iteration.

SGD is ~~more~~ computationally less expensive and converges faster, but is not a smooth descent due to the ~~set~~ arbitrary switching of objective functions, while standard GD is more computationally intensive since it must calculate an average of all of the gradient funcs each iteration, but produces a smoother descent.

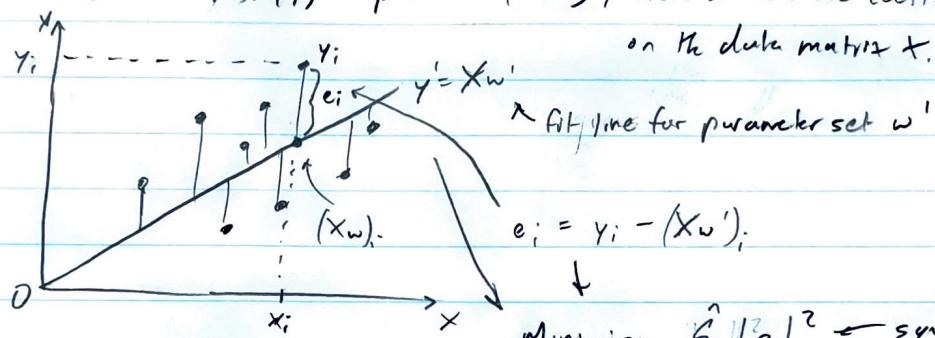
? - Next page

2. a. Linear regression finds the "best" line (through the origin) for the data x_i, y_i for $i = 1 \dots n$ by minimizing the sum of the squared residuals between each data point and the line of best fit:

$$\min_{w \in \mathbb{R}^2} \|y - Xw\|^2 \leftarrow \text{squared error}$$

residuals = "error" = difference between label and prediction

\uparrow
error = label(y) - prediction(Xw), where w are the coefficients / weights on the data matrix X .



$$\text{Minimize } \sum_{i=1}^n |e_i|^2 \leftarrow \text{squared residuals/error strength}$$

The normal equations give w^* :

$$x^T x w^* = x^T y$$

$$\therefore w^* = (x^T x)^{-1} x^T y \quad (\text{Moore penrose Pseudo-inverse})$$

$$\text{where } \hat{y} = Xw^*$$

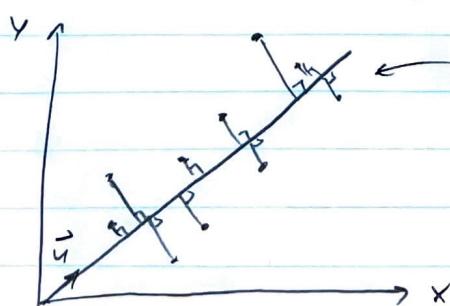
b. $X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} = \begin{bmatrix} x_1 & y_1 \\ \vdots & \vdots \\ x_n & y_n \end{bmatrix}$

i.e. "principal component"

Principal components analysis aims to find the optimal 1-dimensional subspace that maximizes the energy/spread/variance of the data after orthogonal projection onto the subspace

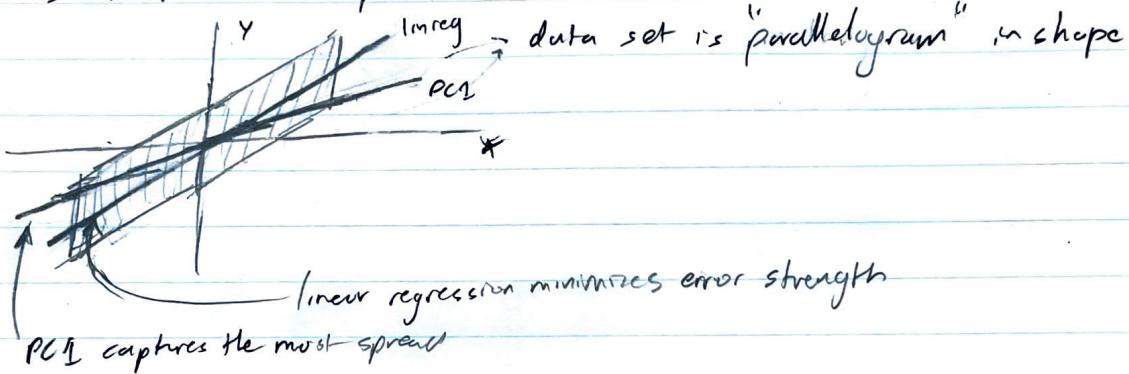
$$\max_{\|u\|_2=1} \sum_{i=1}^n (x_i \cdot u)^2 = \max_{\|u\|_2=1} u^T x^T x u$$

x centered data matrix
 u is direction of 1D subspace / principal component



orthogonal projection of data on to 1D subspace
Find direction \vec{u} that maximizes spread/variance of projected data

(2). c. The lines found from linear regression and PCA are not the same since they are optimizing different quantities — the sum of the errors squared for linear regression ~~is~~ is minimized, while PCA maximizes the variance of the projected data onto a 1D PC/subspace. The line that captures the most spread/variance in the data is not necessarily the same one that ~~not~~ minimizes the ~~the~~ error in ~~the~~ the least squares sense. For example:



d. Linear regression objective function: $L(w) = \|y - Xw\|^2$

PCA objective function: $\epsilon = \sum_{i=1}^n (x_i \cdot u)^2$

Linear regression aims to reduce the error between the data labels and their "predicted" value by optimizing the prediction curve / line of best fit. The error is defined as the difference between the actual values/labels y and the predicted values Xw .

$$e = y - Xw$$

$L(w) = \|e\|^2 \rightarrow$ minimize objective fcn $L(w)$

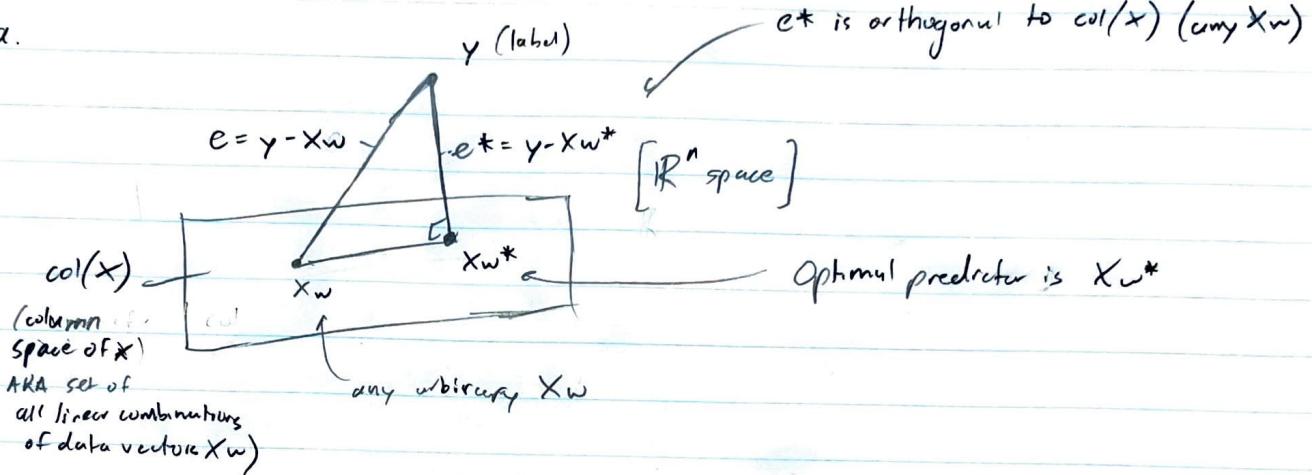
PCA aims to maximize the objective function ϵ , meaning it aims to maximize the variance/spread of the orthogonally-projected data by optimizing the PC line.

$(x_i \cdot u)^2$ represents the distance of the projected point x_i on \vec{u} from the origin.

$$3. e^* = y - Xw^* \in \mathbb{R}^n$$

\uparrow
n x p
data matrix w arbitrary vector
 in \mathbb{R}^p

a.



b. Orthogonality principle states that optimal error vector $e^* = y - Xw^*$ is orthogonal to each row of X^T (each column of X)

$$\text{For orthogonality: } x^T y = 0 = y^T x \rightarrow e^* \cdot X = 0 = X^T e^*$$

$$\begin{aligned} x^T e^* &= 0 \\ x^T (y - Xw^*) &= 0 \\ x^T y - x^T Xw^* &= 0 \rightarrow x^T y = x^T Xw^* \\ I \leftarrow (x^T x)^{-1} x^T Xw^* &= (x^T x)^{-1} x^T y \end{aligned}$$

$$w^* = (x^T x)^{-1} x^T y$$

expression for
optimal least squares
coefficients

where x is data matrix
 y are labels

4. Next page

same as expression from
objective fun and gradient in
lecture

4. $x_i, y_i \in \mathbb{R} \rightarrow (\rho = 1)$

Not centered data

Centroid: $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$, ~~$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$~~ $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$

~~$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$~~

~~$$\begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$~~

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$w = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

append 1's
to center (account
for y-intercept)

$$X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$

data matrix

$$y = Xw + e$$

error

prediction \hat{y}

$$\hat{y} = Xw^*$$

optimal w parameters

$$\hat{y} = Xw^*$$

$$w^* = (X^T X)^{-1} X^T y$$

$$\begin{bmatrix} w_0^* \\ w_1^* \end{bmatrix} = \left(\begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$X^T X = \begin{bmatrix} (1 \times n) & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{bmatrix}$$

$$X^T y = \begin{bmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{bmatrix}$$

~~$$e = y - Xw = 0?$$~~

~~$$\bar{e} = \bar{y} -$$~~

$$\bar{e} = \frac{1}{n} \sum_{i=1}^n e_i = 0?$$

Change

and?

Rearrange to normal equations from $w^* = (X^T X)^{-1} X^T y$

$$X^T X w^* = X^T y$$

See next page

(4). Normal equations : $\mathbf{x}^T \mathbf{x} w^* = \mathbf{x}^T \mathbf{y}$

Expand:

$$\underbrace{\begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{bmatrix}}_{\textcircled{1} \quad \mathbf{x}^T \mathbf{x}} \underbrace{\begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}}_{w^*} \begin{bmatrix} w_0^* \\ w_1^* \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{bmatrix}}_{\textcircled{2} \quad \mathbf{x}^T \mathbf{y}} \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}}_{\mathbf{y}}$$

$$\textcircled{1} \quad \mathbf{x}^T \mathbf{x} = \begin{bmatrix} (1 \times n) & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{bmatrix}$$

$$\textcircled{2} \quad \mathbf{x}^T \mathbf{y} = \begin{bmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{bmatrix}$$

Now:

$$\underbrace{\begin{bmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{bmatrix}}_{\textcircled{1} \quad \mathbf{x}^T \mathbf{x}} \underbrace{\begin{bmatrix} w_0^* \\ w_1^* \end{bmatrix}}_{w^*} = \underbrace{\begin{bmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{bmatrix}}_{\textcircled{2} \quad \mathbf{x}^T \mathbf{y}}$$

From first equation :

$$n w_0^* + \left(\sum_{i=1}^n x_i \right) w_1^* = \sum_{i=1}^n y_i \quad \leftarrow \text{Divide terms by } n$$

$$w_0^* + \underbrace{\left(\frac{1}{n} \sum_{i=1}^n x_i \right)}_{\bar{x}} w_1^* = \underbrace{\left(\frac{1}{n} \sum_{i=1}^n y_i \right)}_{\bar{y}}$$

$$w_0^* + \bar{x} w_1^* = \bar{y}$$

\therefore [optimal least squares line passes through the data centroid (\bar{x}, \bar{y})]

Honor code: On my honor I have neither given nor received any unauthorized aid on this exam.



Midterm Cheat Sheet

1. Linear algebra, optimization

Unsupervised: PCA, k-means clustering, hier. clust.

Supervised: lin reg, least squares, ridge/Lasso reg.

2. Vector space V : α is scalar } ; $\alpha x \in V$, $x+y \in V$ $x, y \in V$

Linear combination: $y = \sum_{m=1}^M \alpha_m x_m$ ✓ vector
 vector \rightarrow i.e. mixing board

Euclidean 2 norm: $\|x\|_2 = \sqrt{\sum_{i=1}^p |x(i)|^2}$ energy $\|x\|_2^2$

ℓ_q -norm: $\|x\|_q = \left(\sum_{i=1}^p |x(i)|^q \right)^{1/q}$

∞ -norm: $\|x\|_\infty = \max_i |x(i)|$ peak value of x

Normalize vector $\|x\|_2 = 1 \rightarrow$ scale by $\frac{1}{\|x\|_2}$

Inner product: $x \cdot y = \langle x, y \rangle = y^T x = \sum_{i=1}^p x(i) y(i)$

Cosine similarity: $\cos(\theta)_{x,y} = \frac{x}{\|x\|_2} \cdot \frac{y}{\|y\|_2} \quad x \cdot y = \|x\|_2 \|y\|_2$

Orthogonality: $x \cdot y = 0 = \langle x, y \rangle$

3. Cauchy-Schwarz inequality: $0 \leq |x \cdot y| \leq \|x\|_2 \|y\|_2$

$x \neq y$ are: most diff or most similar
 when orthogonal or when collinear

Basis: vector space V , lin. ind. & span V

dimension of V is p

$$b_k \cdot b_k = 0 \quad k \neq l$$

Orthogonal basis: $\{b_k\}_{k=1}^p \rightarrow$ elements orthogonal

Orthonormal basis \rightarrow orthogonal basis & normalized $\|b_k\|_2 = 1$

Inverse of matrix: $AA^{-1} = A^{-1}A = I$

$\hookrightarrow B$ contains orthonormal basis, then $B^{-1} = B^T$

and B is orthogonal matrix

$$x = B_a = \sum_{k=1}^p a_k b_k \text{ synthesis; } a = B^T x, a_k = x \cdot b_k$$

Data matrix X $n \times p$ $[x]_{ij} \quad x_i:$ datapoint (vector)

$\Sigma: \mathbb{R}^{p \times p}$ $p \times 1$ column vec.

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix} \quad \text{vector rows}$$

Correlation matrix $X^T X \rightarrow$ symmetric

$$X^T X = \sum_{k=1}^n x_k x_k^T$$

Centering & normalizing: $\mu_i = \frac{1}{n} \sum_{k=1}^n x_k(i)$

$$x_{j[i]} \leftarrow \frac{x_{j[i]} - \mu_i}{\sigma_i}$$

$$\sigma_i = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (x_k(i) - \mu_i)^2}$$

$$3. Av = \lambda v \quad \begin{array}{l} \leftarrow \text{eigenvalues} \\ \uparrow \quad \leftarrow \text{eigenvectors} \\ \text{Must be square} \end{array} \quad A = V \Lambda V^{-1}$$

real symmetric matrix (correlation covariance)
 for example, have real-valued eigenvals & eigenv.

Hermitian transpose A^H
 $A^H A = A A^H = I \rightarrow$ unitary matrix A (extension of orthogonal)

$$\text{SVD: } A = U \Sigma V^H$$

$\uparrow \quad \uparrow \quad \uparrow$
 $n \times p \quad n \times n \quad p \times p$
 unitary unitary \uparrow \uparrow \uparrow
 \uparrow \uparrow \uparrow \uparrow
 w/ eigenvectors w/ eigenvectors w/ singular values (non-reg)
 of $A^H A$ of $A A^H$ real vals

columns of U, V^H
 U, V^H are orthonormal bases

Singular value $\sigma_j = \sqrt{\lambda_j}$ $A^H u_j = \sigma_j u_j$

Frobenius norm $\|A\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^p |\alpha_{i,j}|^2}$

Rank = # of nonzero singular vals.

4. PCA \rightarrow Data centered at 0 + normalized

$$\begin{array}{c} x_i \\ \downarrow \\ O \quad u \\ \uparrow \\ (x_i - u) \cdot u \\ \downarrow \\ \|u\|_2 = 1 \\ \varepsilon = \sum_{i=1}^n (x_i \cdot u)^2 \end{array}$$

$$\text{PCA} \rightarrow \max_{\|u\|_2 = 1} \sum_{i=1}^n (x_i \cdot u)^2 = \|x_c u\|^2$$

maximize energy/spread/
 variance of data after orthogonal projection onto

$$x^T x = V \Lambda V^T \rightarrow \max_{\|u\|_2 = 1} u^T V \Lambda V^T u$$

1st subproj

covariance matrix Σ Σ is first eigenvector v_1 of $x^T x$

$$\text{energy is eigenvalue} \Rightarrow \lambda_1 = \sum_{i=1}^n (x_i \cdot v_1)^2$$

2nd PC is orthogonal to first $\&$ is 2nd leading eigenvector

PCA via SVD: $x = U \Sigma V^T$

$$\begin{array}{c} \uparrow \quad \uparrow \\ \sigma_k^{-2} \text{ is energy in PCs} \quad \text{right singular vectors} \\ \text{columns of } V \text{ are PCs} \end{array}$$

$$P = X V = U \Sigma$$

5. Objective (lin.) fun. L , w_0 is minimizer/max...

Convex fun.: for any x_1, x_2 in f and $0 < t < 1$

$$f(tx_1 + (1-t)x_2) \leq t f(x_1) + (1-t) f(x_2)$$

Concave set: $t x_1 + (1-t)x_2 \in S$

$$\text{Local min } \|w - w^*\| \leq S$$

constraint

$$\text{Lagrange multiplier } \lambda: \min_{w \in \mathbb{R}^p, \lambda \in \mathbb{R}} L(w) + \lambda g(w)$$

$$\text{Gradient: } \nabla L(w) = \left[\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \dots, \frac{\partial L}{\partial w_p} \right]^T$$

$$= w^+ - \mu + \sum_{t=1}^T \nabla L_t(w)$$

↓
step size / learn rate

$$(5) \text{ Gradient descent: } w^{t+1} = w^+ - \mu + \nabla L(w^+)$$

↳ Issues w saddle and not smooth

Stochastic GD: $L(w) = \sum_{t=1}^T L_t(w)$ ← sum of smaller obj. functions
 ↳ computationally less expensive; slow convergent & non monotonic

$$w^{t+1} = w^+ - \mu + \sum_{t=1}^T \nabla L(w)$$

↳ Epoch = all t used
 ↳ set of gradients picked each iteration (sequence or random)

6. Cluster → homogeneous groups; unsupervised

K-means: $\min_{C_1, \dots, C_K} \sum_{k=1}^K \sum_{i=1}^n \|x_i - C_k\|_2^2$ ← within cluster variation minimize

↳ $L(C_k)$ w/ 2 norm = $\frac{1}{|C_k|} \sum_{i \in C_k} \|x_i - \bar{x}_k\|^2$

↳ computationally inefficient

finds local min

- 1) Random cluster assignment,
- 2) Centroid/ $\frac{1}{|C_k|} \sum_{i \in C_k} x_i$
- 3) Assign to closest centroid

Different "distance" metrics
 i.e. Euclidean $\sqrt{\sum (x_i - y_i)^2}$

- 1) $\|\cdot\|_2$ 2-norm
- 2) $\|\cdot\|_2^2$ squared 2
- 3) $\|\cdot\|_1$ 1-norm
- 4) $\|\cdot\|_\infty$ max norm (max dist)
- 5) Cosine distance

Hierarchical clustering → tree-based taxonomy

↳ Agglomerative (bottom-up) ← Divisive (top-down)

↳ Similarity function → dendrogram tree

8. $x_i \rightarrow$ data point/vector $y_i \rightarrow$ label

$$y_i = w_1 + w_2 x_i + e_i \leftarrow \text{minimize error strength } \sum_{i=1}^n \|e_i\|^2$$

$y = Xw + e$ (least squares)

$$\min_w \|y - Xw\|^2 \rightarrow L(w) = \|y - Xw\|^2$$

Objective fn

$$\nabla L(w) = 0 \text{ @ minimizer } w^*$$

$$\nabla L(w) = \left[\frac{\partial L(w)}{\partial w_1}, \frac{\partial L(w)}{\partial w_2} \right]^T$$

$$L(w) = \|y - Xw\|^2 = (y - Xw)^T(y - Xw)$$

$$y^T \lambda w = w^T X^T y \leftarrow y^T y - y^T Xw - w^T X^T y + w^T X^T Xw$$

$$= y^T y - 2w^T X^T y + w^T X^T Xw$$

$$\nabla_2 z^T a = a, \quad \nabla_2 z^T A_2 = 2A_2$$

$$\nabla L(w) = -2X^T y + 2X^T Xw = 0 \rightarrow X^T y = X^T Xw$$

Normal eqns: $X^T y = X^T Xw^*$

$$w^* = (X^T X)^{-1} X^T y$$

$\underbrace{w^*}_{= X^+ \text{ Moore Penrose}}$

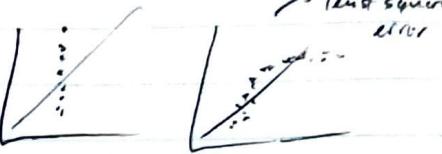
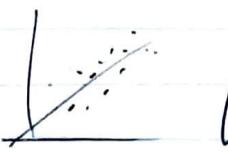
$$X = U\Sigma V^T, \quad X^+ = V\Sigma^+ U^T$$

↳ reciprocal of non-zero entries of Σ^+

$$\text{Quadratic} \quad X^T X w + b$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$$

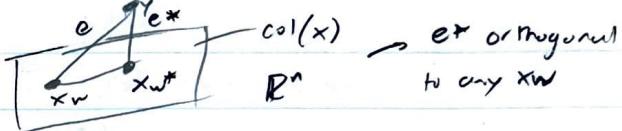
$$8. w^* = (X^T X)^{-1} X^T y \quad \hat{y} = Xw^*$$



→ least squares error

Orthogonality principle $X^T e^* = 0$

↳ optimal error vector e^* is orthogonal to each row of X^T and ∴ each column of X



Include column of 1s for y-intercept term w_0

Center data by subtracting mean of y and each column of x

$$I^T y = 0, I^T X = 0^T$$

→ w is identifiable and w's solve numerically

$(X^T X)^{-1}$ does not exist → $X^T X$ is singular or ill-conditioned

1) $n \geq p$ (full, numeric x) but $\text{rank}(X) < p$

↳ 2+ columns of x (and $X^T X$) are lin. dependent

2) $n < p$ (short, wide x) → overfitting; many params w_j

↳ then data points (x_i, y_i) model fit to noise

↳ axes are eigenvectors

w_1, w_2 eccentricity is from condition #

$$w_{1,2} = w^*$$

student x
 variance
 bins

Regularization w/ constraint: $\min_{w \in \mathbb{R}^p, \lambda \in \mathbb{R}} \|y - Xw\|^2 + \lambda \|w\|_2^2$

Ridge $X^T y = (X^T X + \lambda I) w_{\text{ridge}} \rightarrow w_{\text{ridge}} = V(\Sigma + \lambda \Sigma)^{-1} U^T$
 always well-conditioned $\|w\|_2^2 = C$

Lasso → 1 norm penalty $\|w\|_1 = C$

w_{lasso} is sparse, "variable selection"

Elastic net → combine 2 penalties

Need optimization alg. for lasso