Lecture 2 -- Statistical Learning

Tuesday, September 3, 2024 3:29 PM

For e-vals/vecs
$$X = QDQ^T$$

 $Xv_i = \lambda_i v_i$ $XQ = QD$

For singular vals/ves -
$$\frac{1}{2}$$
 - $\frac{1}{2}$ $\frac{1}{2}$

$$\frac{8p}{X}$$
 = $\begin{pmatrix} 3 & 2 \\ 2 & 3 \end{pmatrix}$ Let's get the SVD.

$$x^{T}x = \begin{bmatrix} 17 & 8 \\ 8 & 17 \end{bmatrix} \qquad \begin{bmatrix} a & b \\ b & a \end{bmatrix}$$

notes Page 1

$$\sqrt{=1/2} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}; D = \begin{bmatrix} 5 & 0 \\ 3 & 3 \\ 0 & 0 \end{bmatrix}$$

$$3x2$$

Know:
$$Xv_i = 6iU_i$$

then $U_i = Xv_i$
 $U_1 = \frac{Xv_1}{6_i} = \frac{1}{5} \left(\frac{32}{2 \cdot 3}\right) \left(\frac{1}{1}\right) \sqrt{52} \left[\frac{5}{5}\sqrt{2}\right] \sqrt{52} \left[\frac{5}{5}\sqrt{2}\right] \sqrt{52} \left[\frac{32}{2 \cdot 2}\right] \left(\frac{1}{1}\right) \sqrt{52} \left[\frac{5}{5}\sqrt{2}\right] \sqrt{52} \left[\frac{32}{2 \cdot 2}\right] \left(\frac{1}{1}\right) \sqrt{52} \left[\frac{5}{5}\sqrt{2}\right] \sqrt{52} \left[\frac{5}{5}\sqrt{2}\right] \sqrt{52} \left[\frac{32}{2 \cdot 2}\right] \left(\frac{1}{1}\right) \sqrt{52} \left[\frac{5}{5}\sqrt{2}\right] \sqrt{52} \left[\frac{5}{5$

notes Page

$$u_2 = \frac{\chi v_2}{6_2} = \frac{1}{3} \begin{bmatrix} \frac{3}{2} & \frac{2}{3} \\ \frac{2}{2} & \frac{2}{3} \end{bmatrix} \begin{pmatrix} 1 \\ -1 \end{pmatrix} / \sqrt{2} \begin{bmatrix} \frac{1}{3} & \frac{1}{2} \\ \frac{1}{3} & \frac{1}{3} \end{bmatrix}$$

Uz is whatever is orthog. to U, and Uz.

$$U_3 = \frac{1}{\sqrt{24\frac{1}{4}}} \left[-\frac{1}{2} \right]$$

Verify at home: X=UDV.

Statistical Machine Learning
Building/analyzing delta analysis procs
Using probabilistic Hinking

(where helpful)

Broadly two main cats of ML probs

	Moderly Two Main Can of 1000
_	> (I) supervised learning
	supervised means we have "examples" to train the ML method
	to train the ML weethood
	idea: want to predict some Y from X and we have examples
	(trains data): (Zn, yn) n=1
	W/in supervised two main types
	(1) <u>regression</u> problems: Y is
	a continuas /real-valved var.
	2) classification problems:
	<u> </u>
	y is discrete valued/finite
	is discrete valued/finite (categorical)
	Ex. regression: YER Y

- predict stock market perf. from economic indicators predicting adult height from chidhood nutrition

Ex. classication! Y is categorical

- predict if individual will

VE \$0,13 default on loan given

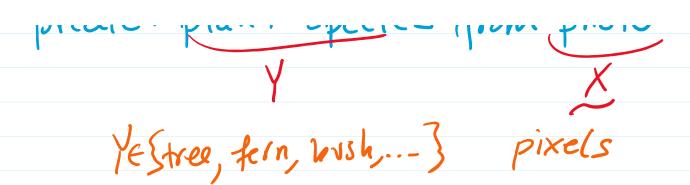
credit score

binary

classification

problem

- predict plant species from photo



II) Unsupervised Problem

No clear prediction problem.

No clear distinction bottom X and Y

(just have some data)

Goal: learn/summarize important trends/rels/info in my data.

Math setup for supervised learning

We have some var Y we want to predict from

predict from

$$X = (X_1, X_2, ..., X_p)$$
 $X = (X_1, X_2, ..., X_p)$
 $X = (X_1, X_2, ..., X_p)$

We're going to come up y a function f so that

 $f = f(x) = f(x_1, x_2, ..., x_p)$

Course goals: talk about 2 things

(1) Methods: how do we construct f ?

(2) Evaluation: how do we determine if f is good?

For super problems we assume we have training data to construct f :

 $f = f(x_1, x_2, ..., x_p)$

notes Page 7

 $(\chi_n, y_n)_{n=1}$ N = Size of training deta where $\chi_n \in \mathbb{R}^p$ $Y_n \in \mathbb{R}^p$ $Y_n \in \{C_1, C_2, ..., C_K\}$ for class.

Linear Regression (regression prohlans)

why? (1) classic method-vvell studied

(2) simple (good)

(3) powerful

(4) basis for more complex wellood

Setup: For lin. reg. assume we have some input vars (design):

$$X = (X_1, ..., X_p) \in \mathbb{R}^p$$

ord assoc. Coef.

then lin. res. assures a model of the form

$$Y = f(X) = X \beta$$

$$= \sum_{j=1}^{P} \beta_j X_j'$$

to learn f we will "learn" some good valves for B, call them B and then form f as

$$\hat{f}(\chi) = \chi \hat{\beta} = \sum_{j=1}^{p} \hat{\beta}_{j} \chi_{j}.$$

Where is the intercept? Why not $f(X) = \hat{\beta}_0 + \sum_{j=1}^p \hat{\beta}_j X_j$ Can hide the intercept in the dasign In this course we use X in two ways (1) vars. we measure
(2) the input to our ML algo. this more
(design) Need not be the same, can do only pre-proc. on (1) to get (2) vais I measure: X1, X2, ..., Xp inpot to ML: X1/ lo(X2), Sin(X3)...

Can always get intercept by setting
$$X_1 = 1$$
 always,

1.l. $X = (1, other vars)$

How do I determine good values for B?

Many ways to determine if Y=f(x)

Simplest way: Least-Squares regression Choose & as valve that minimizes the squared training error (loss)

$$L(b) = \sum_{n=1}^{N} (y_n - x_n \beta)^2$$
trefarget predicted in training val.