## ML Assignment

```
Singular value decomposition (SVD) factorizes a mixin matrix X as X = U \in V^T, where U \in E^{m \times n} and U^TU = UU^T = T, E \in E^{m \times n} Contains non-increasing non negative values along its diagonal and zeros essewhere, and V \in R^{n \times n} and V^TV = VVT = T. Given the SVD of a matrix X = U \in V^T, what is the eigendcomposition of XXT? (You should define an appropriate square matrix Q and diagonal matrix A such that XXT = Q \times Q^{-1}).
```

Here, Given that,  $X = U \leq V^T$  where  $U \in P^{m \times h}$  and  $U^T U = U U^T = T$ ,  $\xi \in P^{m \times h}$ also,  $V^T V = V V^T = T$ 

where, v is caued eigen vector and  $\xi^2$  is caued eigen valuel.

requ

70

T.

аt

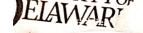
'n

Positive (semi) - Definite Madrices (10 pts each past) let A be a real, symmetric dxd matrin. we say A is the Semi-definite (PSD) if, for all nepd, xTAXZO. We say A is the definite (PD) if, for all x 10, x TAX >0. we with A 20 when A is PSD, and A > 0 when A is PD. The spectral theorem says that ewy real symmetric matrix A can be expressed A = UNUT, where vis a dxd matrix such that UUT=UTU=T (aned an Orthogonal matrix), and  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$ . Multiplying on the night by U and  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$ . The denote the its alumn we see that  $\Delta v = U \otimes \Lambda$ . If we tell u; denote the its alumn expression reveals q u, we have  $\Delta u$ ;  $\Delta u$ ; but each i. This expression reveals that the is one eigenvalues of A, and the corresponding columns ui are eigenvectors associated to di using the spectral decomposition, show that a) A 16 POD iff A; 20 for each i. Hint: UAUT = 2 = 2; u; u; T (you do not need to place this But please try to understand why this holds. Univen, A is pso if for all xerd, XTAXZO. here, and A is PD If for all, x = 0, XTAX>0. so, when, A is pad, we write A > 0.

when A is po, we write A > 0. a) A is pso iff a; ≥0 for each i. proof: By spectral theorem, A = UNUT - Bheie UNIO a dxd matrix Such that UUT = UTU = I and n = diag(2). musplying equia on both sides by U. on  $AU = U \wedge U^{T}U$   $AU = U \wedge U^{T}U$   $U^{T}U = I$ Athra we need to show, lizo if Alis PSD, PED properties holds that: for all xerd, xTAXZO Hence, Mulipiging equilby by UT on both, UTAU = UTUA [ ... UTU=I]

tions of an is with a

Where, A 20 [every element of 1 i.e 2; 20] SO, we get UTAUZO So,  $\lambda; \geq 0$  Hence proved Second past To prove that A is poor if A; ≥ 0 Given that by operal the over, A = UNUT MULLIPLY BY YTON BOTH Sides Again, musipy by y to get it in in the form q: XTAX  $y^TAy = y^TUAU^TY$  — (1) NOW, Let P = UTY Then, Eqn (1) is Let's assume  $p = \begin{cases} p_1 \\ p_2 \end{cases}$ and,  $pT = \begin{cases} p_1, p_2, \dots, p_M \end{cases}$ also, we have  $\Lambda = \begin{cases} 0 & \lambda_1 & 0 \\ 0 & \lambda_2 \end{cases}$ PTAP = [P1, P2...Pa] [ \lambda\_L O O O ] [P2] Fry P2. Pd ] A1 P1
A2 P2
AdPd = 12 P12 + 22 P22 + 20 Pd2  $= \sum_{i=1}^{d} \lambda_i P_i^2$ (a) is:  $y^T A y = \sum_{i=1}^{d} \lambda_i P_i^2$ hence, ey" (2) is:



Manimum Likelihood Estimation (15 pts). Consider a random vasiable x (possibly a vector) whose distribution Lidensity function or mass function) belongs to a parametric family. The density or mass function may be written for f(x;0), where 0 is larted the parameter, and can be either a salar or vector. For example, in the Gaussian family, o can be a two dimensional vector consisting of the mean and vasiance. Suppose the pasametric family is known, but the value of the pasameter is unknown. It is often of interest to estimate this pasameter from observations q X.

Marimum likelihood estimation is one of the most important

pasameter estimation techniques. Let X1,... Xn be i.i.d (independent 1 identically distributed) random vasiables distributed according to  $f(x; \theta)$ . By independence, the joint distribution of the

observations is the product.

appear regions 20 autam) ite

ff f (X; 30) — (1) (6. 11) (6. 11). viewed as a function q 0, this quantity is caused the likelihood q 0. It is often more convenient to work with the log-likelihood is logf(Xi; 0) — (2)

A manimum dikelihood estimate (MLE) of 0 is any parameter Deagmano ( Log f(x; i) -(3))

"as a max" denotes the set of an values achieving the maximum. If there is a unique maximizes, it is cauld the maximum dikelihood estimate. Let x1,... Xn be iid poisson vasiables with intensity pasameter 2- Determine the maximum likelihood estimator qu'd.

moder, and it is about signing the first of

Sola is to find out the maximum likelihood estimator of the posameter of poisson distribution.

poisson's distribution is given by:

$$f(x) = \frac{e^{-\lambda} x^{\chi}}{x!} \quad \text{where, } x = 0, 1, 2, \cdots$$

Goal: To find the maximum likelihood estimator of parameter & from the observations:

ai are independent identically distributed raines

Each q xi is a realization of random variable that has poission distribution.

Now, log likelihand function is given by:

$$L(\lambda) = L_{n} \pi^{f(\alpha; \lambda)}$$

$$= \underbrace{\underbrace{2}_{i=1}}_{\lambda_i} \underbrace{L_0 e^{-\lambda} \lambda^{\chi_i}}_{\chi_i!}$$

$$= \sum_{i=1}^{2} \frac{1}{\ln e^{-\lambda}} + \sum_{i=1}^{2} x_i \ln \lambda - \sum_{i=1}^{2} \ln x_i!$$

arg max UZ)

$$\frac{\partial}{\partial \lambda} L(\lambda) = 0$$

$$\frac{\partial}{\partial \lambda} \left\{ \frac{\partial}{\partial \lambda} (-\lambda) + \frac{\partial}{\partial \lambda} \chi; Ln\lambda - \frac{\partial}{\partial \lambda} Ln \chi; \right\} = 0$$

$$\sigma_{1}, -n + \frac{1}{\lambda} \stackrel{?}{\underset{i=1}{\xi}} \chi_{i} = 0$$

$$\therefore \hat{\lambda} = \frac{1}{n} \stackrel{?}{=} x; \quad (which is sample mean)$$

o, The maximum likelihood estimator of the pasameter 2 for a poisson distributed random vas; able 2 is given by the sample mean of no observations.

To this problem you will proke some a properties of unconstrained optimization problems. For the next two pasts, the following 4) unconstrained optimization fact will be hapful. A twice continuously differentiable fund too  $f(x) = f(y) + (\nabla f(y), \pi - y) + \frac{1}{2} (\pi - y, \nabla^2 f(y)(\pi - y)) + o(\ln \pi - y \|^2)(y)$ admits the quadratic expansion. where old denotes a function satisfying tim old =0, as well f(x) = f(y) + ( \f(y), x-y) + \frac{1}{2}(x-y), \frac{7}{2}f(y+t(x-y))(x-y)) (s) as the expansion a) Show that if fis twice continuously differentiable and xx is a local minimizes, then  $\nabla^2 f(a^*) \ge 0$ , i.e the nession of f is the Semi-definite at the local minimizes 2. hint based on the equal we can write,  $f(x^* + ty) = f(x^*) + (\nabla f(x^*), ty) + \frac{t^2}{2}(y, \nabla^2 f(x^*), y) + o(t^2)$ ) show that if f is twice continuously differentiable, then f is conven if a only if the nessian 72 f(x) is positive semi-definite for all ne Rd. Hint = the proof has two steps . Step 1 . Assume that I is conex. Then prove that 72f(x) is psd- In this step, you should stast with the following, fatty) = f(x) + (vf(x), ty) For the left hand side, we can use equation (4). Step 2. Assumi that  $\nabla^2 f(\eta)$  is psd. Then show that fis conex. To do that ) consider the function  $f(x) = 1/2 x^T A x + b^T x + c$ , where A is a symmetric dxd matrix. Derive the Hessian q f. under

what conditions on A is f convex? Strictly convex?

Since, y'H(x+) y z o ond x(x+, 1y) -se 1 - 50

B K 1 110 - 1 0 x (\*\*) + (y1 1 \* x ) ).

i die flice entry in me) .

ideal to reach the Portion Local

mains rada a si sidi

The Country of

```
b) -> Proof:
      that, f is twice continuously different able.
  > Here,
      we make the use of 2nd order taylor series expansion of (4)
f(x^{*}+ty)=f(x^{*})+(xf(x^{*}),ty)+\frac{t^{2}}{2}(y), x^{2}f(x^{*})+y)+o(t^{2})
        f(x) = f(y) + (\nabla f(y), x - y) + \frac{1}{2}(x - y), \nabla^2 f(y) (x - y)) + o(||x - y||^2)
             If x + is a local minimizes then,
        The to second order Taylor series expansion of that
            a given point x + eIR is given by:
           f(n) = f(x^*) + \nabla f(x^*)^{T} (x - x^*) + \frac{1}{2} (x - x^*)^{T} \nabla^{2} f(x^{2}) (x - x^{2})
                             +0(112-x+1/2) - (a)
           where, Lim O(||x-x*||2) =0.
       Now, Vfat) = 0 because If x* is a local minimum qf,
            then of(xx)=0
        Given, yer a t>0
             let 2 = x* + ty, put this into eqn (a)
           we get 0 \le \frac{f(x^* + ty) - f(x^*)}{12} = \frac{1}{2} y^T \nabla^2 f(x^*) d + o(t^2)}{t^2}
           as \nabla f(x^*) = 0, Take the limit as t \Rightarrow 0, we get
                      0 5 4 72 f(x+) y
           y way chosen arbitrasily, 72f(x*) is positive se
                         definite
     which means that the nession of fis PSD at the
              local minimizes xx
                           Photed
```

```
4. b) -> Proof:
  to prove f is conven if a only if Hessian of 2 for is ASD for
     all xerd.
   NOW, suppose that f is convex and next then by
        f(x+ty) \ge f(x) + (\forall f(x), ty)
    on floctty) > f(n) + + \ f(n) Ty - (a)
     Then, Replacing the L.H.S of the inequality (9) with its
      second order Taylor expansion which gives:
       f(n) + t of(n) y + + +2 y 1 to v2 f(n) y + o(+2) > f(n) +
                           t of couty
   or, \frac{1}{2}y^{t}\nabla^{2}f(n)y + \frac{O(t^{2})}{t^{2}} \geq 0
        to yield,
    As aywas arbitrory, 72 fca) 18 positive semi definite.
      Conversely, if xyerd, then by mean value Theorem
       there is a 6 E (0,1) such that
         f(n) = f(y) + \f(y)^T(x-y) + \frac{1}{2}(n-y)^T f(ye)(n-y)
               which is eq (5)
       where, 46 = 6x + (1-6)4
              f(n) z f(y) + vf(y) T(n-y)
       As, 72 f(ye) is PSD, f is conux. by the
     f(n) zf(y) + \f(y)^T(n-y) for all n, y \end{all}
     according to the ex":
        poul
```

yes nese, binen that,  $f(n) = \frac{1}{2}x^TAx + b^Tx + C$ and A is symmetric dxd matrix To derive nession 9 f. let us consider a magin n as. VA(x) = \[ \frac{2}{2} \times we have a relation in matrix notation, Vf(x) = ATX+AX=(AT+A)x and,  $\nabla f(x) = (A^T + A) x$ and, if A is symmetric then, AT=A so, we have, (AT+A) = (A+A) = 2A Hence, - (4)  $\nabla f(n) = 2Ax - (1)$ Then, consider the above according egn,  $f(n) = \frac{1}{2} n^T A x + b^T x + C$ Using the knuts from equ(1) we knu,  $\nabla f(x) = \frac{1}{2} (AT + A) x + b \qquad \left[ \frac{d(x^T a)}{dx} = \frac{d(a^T x)}{dx} = a^T \right]$ or,  $\nabla f(\alpha) = \frac{2}{An+b}$  [... A is symmetric]  $\nabla^2 f(x) = H(x) = A$ Again, To find the conditions under which is con un, a strictly convex The function  $f(x) := \frac{1}{2} x^T A x + b^T x + c$  is a convex function if and only if A is positive semi The function  $f(n) = \pm x^T A x + b^T x + c$  is strictly conux if A is positive definite and, only if A > 0 % definite.