

DSA3102 Lecture 10 (Reading Sections 7.1 & 7.5)

- Convex Sets / Functions / Optimization Problems
- Duality³ (Lagrangian, Dual Opt. Problem, KKT)
- Statistical Applications
- Optimization Algorithms (Unconstrained, Equality, Inequality).
 Gradient / Steepest Descent Newton's method.

Family of probability distributions on \mathbb{R}^m indexed by

parameter vector $x \in \mathbb{R}^n$, $\{p_x(\cdot)\}_{x \in \mathbb{R}^n}$

argument is y here

For each $x \in \mathbb{R}^n$, $p_x(\cdot)$ is a probability distribution.

Continuous case: $p_x(y) \geq 0 \quad \forall y \in \mathbb{R}^m, \int_{\mathbb{R}^m} p_x(y) dy = 1 \quad \forall x \in \mathbb{R}^n$ hold for every parameter

Discrete case: $p_x(y) \geq 0 \quad \forall y \in \mathbb{R}^m, \sum_y p_x(y) = 1.$

More convenient to work with log-likelihood

make dependence
on y explicit

$$l_y(x) \triangleq \log p_x(y)$$

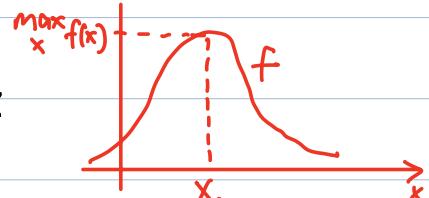
y : observations / measurements
 x : unknown parameters.

Maximum likelihood estimation (MLE)

$$\hat{x}_{ml}(y) = \underset{x \in \mathbb{R}^n}{\operatorname{argmax}} p_x(y) = \underset{x \in \mathbb{R}^n}{\operatorname{argmax}} \log p_x(y)$$

opt prob?

$$= \underset{x \in \mathbb{R}^n}{\operatorname{argmax}} l_y(x) \leftarrow \text{log likelihood of } x$$



If we have prior information about x ,

say $x \in C \subseteq \mathbb{R}^n$, this can be incorporated as

constraints in (*) as follows

make constraint explicit in the opt prob
itself

$$\hat{x}_{ml}(y) = \underset{x \in C}{\operatorname{argmax}} l_y(x)$$

$\underset{x}{\operatorname{argmax}} f(x)$
(argument that maximizes $f(x)$)

Examples: Linear measurements in i.i.d. noise.

$$y_i = a_i^T x + v_i, \quad i=1, \dots, m$$

$x \in \mathbb{R}^n$: unknown vector of parameters

$y_i \in \mathbb{R}$: observations (known) attributes

$a_i \in \mathbb{R}^n$: feature vectors (known) v_i 's independent

$v_i \in \mathbb{R}$: are measurement errors or noise with distribution $p_v(\cdot)$.

(y_1, \dots, y_n) probabilities

$$\text{Likelihood: } p_x(y) = \prod_{i=1}^m p_v(y_i - a_i^T x)$$

randomness is
 v_i express as that

likelihood of observing measurement vector y
= likelihood of observing that measurement errors all $y_i - a_i^T x$

$$\text{Log-likelihood: } l_y(x) = \log p_x(y) = \log \left(\prod_{i=1}^m p_v(y_i - a_i^T x) \right)$$

$$= \sum_{i=1}^m \log p_v(y_i - a_i^T x)$$

ML estimate is the optimal solution for

f need not
be convex

$$\max_{x \in \mathbb{R}^n} \sum_{i=1}^m \log p_v(y_i - a_i^T x)$$

Def: A function $f: \mathbb{R}^n \rightarrow (0, \infty) = \mathbb{R}_{++}$ is log-convex if

$g: \mathbb{R}^n \rightarrow \mathbb{R}$ defined as $g(x) = \log f(x)$ is a convex function.

Properties: f_1, f_2 are log-convex, $f = f_1^\alpha f_2^{1-\alpha}$ is log-convex for any $\alpha \in [0, 1]$.

convex combination

$$\begin{aligned} \text{Pf: } g(x) &= \log f(x) = \log ((f_1(x))^\alpha (f_2(x))^{1-\alpha}) \\ &= \underbrace{\alpha \log f_1(x)}_{\text{convex}} + (1-\alpha) \underbrace{\log f_2(x)}_{\text{convex}} \rightarrow \text{convex.} \end{aligned}$$

Def: A function $f: \mathbb{R}^n \rightarrow (0, \infty) = \mathbb{R}_{++}$ is log-concave if

$g: \mathbb{R}^n \rightarrow \mathbb{R}$ defined as $g(x) = \log f(x)$ is a concave function.
so that we can take log

Folklore: Most densities/distributions are log-concave.

(most of the
time it's true)

Example: Suppose v_i 's are Gaussian with zero mean

and variance σ^2 ($E[v_i] = 0$, $\text{Var}(v_i) = \sigma^2$), the density

$$p_v(v_i) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{v_i^2}{2\sigma^2}\right), \quad v_i \in \mathbb{R}.$$

Probability
density function

density
of Gaussian

log concave

arg max (ve of quadratic)

$$\text{MLE: } \hat{x}_m = \arg \max_{x \in \mathbb{R}^n} \sum_{i=1}^m \log p_v(y_i - a_i^T x)$$

$$v_i = y_i - a_i^T x$$

$$= \arg \max_{x \in \mathbb{R}^n} \sum_{i=1}^m \log \left(\frac{1}{\sqrt{2\pi} \sigma} \exp \left(- \frac{(y_i - a_i^T x)^2}{2\sigma^2} \right) \right)$$

$$= \arg \max_{x \in \mathbb{R}^n} \left(-\frac{m}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|Ax - y\|^2 \right)$$

does not involve x (not important)

A is the matrix with rows a_1^T, \dots, a_m^T .

ML estimate of x : $\hat{x}_{ml} = \arg \min_{x \in \mathbb{R}^n} \|Ax - y\|^2$ complex opt prob

least squares problem.

Laplacian noise

$$p_v(v_i) = \frac{1}{2a} e^{-|v_i|/a}, \quad v_i \in \mathbb{R}$$

non differentiable
where

$$\frac{1}{2a}$$

or $f(v_i)$

symmetric

decreasing
linear in v_i

v_i

log-concave

$$\text{MLE: } \hat{x}_{ml} = \arg \max_{x \in \mathbb{R}^n} \sum_{i=1}^m \log p_v(y_i - a_i^T x)$$

$$-|y_i - a_i^T x|/a$$

minus l_1 norm

concave as l_1 norm is convex

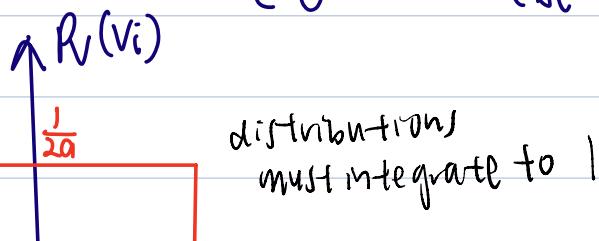
$$= \underset{x \in \mathbb{R}^n}{\operatorname{argmax}} \sum_{i=1}^n \log \left(\frac{1}{2a} e^{-|y_i - a_i^T x|/a} \right)$$

$$\hat{x}_{\text{ML}} = \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} \|Ax - y\|_1$$

When the noise is Laplacian, the MLE is an ℓ_1 -norm approximation problem.

$$\|z\|_1 = \sum_{i=1}^m |z_i|,$$

Uniform noise uniform dist $p_V(v_i) = \begin{cases} \frac{1}{2a}, & v_i \in [-a, a] \\ 0, & \text{else} \end{cases}$



holds true for any dist

MLE: $\hat{x}_{\text{ML}} = \underset{x \in \mathbb{R}^n}{\operatorname{argmax}} \sum_{i=1}^m \log p_V(y_i - a_i^T x)$

either constant or 0
 $\log 0 \rightarrow -\infty$, never achieve the arg max)

ML estimate: \hat{x}_{ML} is any vector x that satisfies $\|Ax - y\|_\infty \leq a$.

feasibility problem

$$\|z\|_\infty = \max_{1 \leq i \leq m} |z_i|$$

If form $Ax - y$, every component should have abs value not more than a

if this satisfied, probability is the same

Counting Problems with Poisson distribution

$y \in \{0, 1, \dots\}$ observations. — Observations take values on a non-negative integer

$$Pr(y=k) = \frac{e^{-\mu} \mu^k}{k!}, \quad k=0, 1, 2, 3, \dots \quad E[y] = \mu$$

y : represent the count or number of events of a Poisson process over a fixed period of time.

Model the mean $\mu = E[y] = a^T u + b \quad u \in \mathbb{R}^n$

u : vector of explanatory variables.

$$\mu = a^T u_i + b.$$

$a \in \mathbb{R}^n, b \in \mathbb{R}$: unknown parameters.

pairs of explanatory variables and measurements

Given # of observations $(u_i, y_i), i=1, \dots, m$.

Aim: Find parameters $a \in \mathbb{R}^n, b \in \mathbb{R}$ from these data.

put μ into Poisson formula

$$\text{Likelihood function: } L_{u,y}(a, b) = \prod_{i=1}^m \frac{e^{-(a^T u_i + b)} (a^T u_i + b)^{y_i}}{y_i!}$$

Log-likelihood function:

$$\lambda_{u,y}(a, b) = \sum_{i=1}^m [y_i \log(a^T u_i + b) - (a^T u_i + b) - \log(y_i!)]$$

$$\begin{bmatrix} a \\ b \end{bmatrix}^T \begin{bmatrix} u_i \\ 1 \end{bmatrix}$$

concave

ML estimate:

$$(\hat{a}_{ML}, \hat{b}_{ML}) = \underset{a, b}{\operatorname{argmax}} \sum_{i=1}^m [y_i \log(a^T u_i + b) - (a^T u_i + b)]$$

log (mean function)
is concave

This is a convex optimization problem.
taking max & the objective f^* is concave.

probability mass function
of poisson []
log-concav

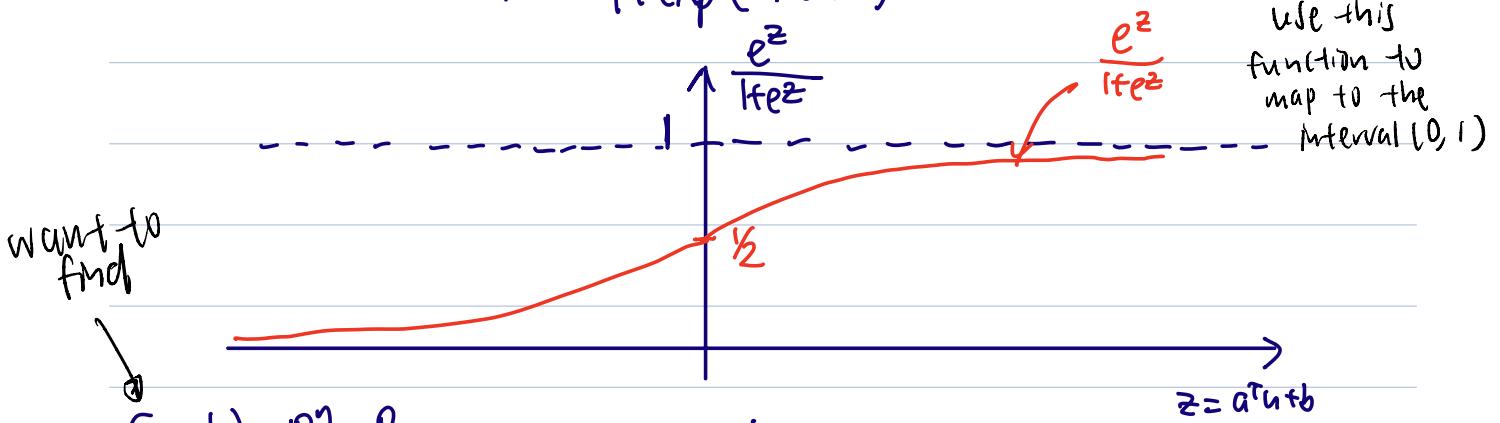
Logistic Regression.

Consider a random variable $y \in \{0, 1\}$

$$\Pr(y=1) = p, \quad \Pr(y=0) = 1-p.$$

$p \in [0, 1]$ depends on a vector of explanatory variables $u \in \mathbb{R}^n$
real number

Logistic model $p = \frac{\exp(a^T u + b)}{1 + \exp(a^T u + b)} \in (0, 1)$



explanatory variables

Data: $u_1, \dots, u_m \in \mathbb{R}^n$ and corresponding outcomes $y_1, \dots, y_m \in \{0, 1\}$.

Given data $\{(u_i, y_i)\}_{i=1}^m$, find (a, b) . (Want a good estimate for the model parameter)

This task is called logistic regression.

Reorder the data s.t. $y_1 = \dots = y_s = 1, y_{s+1} = \dots = y_n = 0$

Likelihood: $L_{u,y}(a, b) = \prod_{i=1}^q p_i \prod_{i=q+1}^m (1-p_i)$

 $= \left(\prod_{i=1}^q \frac{\exp(a^\top u_i + b)}{1 + \exp(a^\top u_i + b)} \right) \left[\prod_{i=q+1}^m \frac{1}{1 + \exp(a^\top u_i + b)} \right]$

$y_i = 1$ $y_i = 0$

Log-likelihood:

 $\ell(a, b) = \log \left[\left(\prod_{i=1}^q \frac{\exp(a^\top u_i + b)}{1 + \exp(a^\top u_i + b)} \right) \left[\prod_{i=q+1}^m \frac{1}{1 + \exp(a^\top u_i + b)} \right] \right]$
 $= \sum_{i=1}^q \log \left(\frac{\exp(a^\top u_i + b)}{1 + \exp(a^\top u_i + b)} \right) + \sum_{i=q+1}^m \log \left(\frac{1}{1 + \exp(a^\top u_i + b)} \right)$
 $= \sum_{i=1}^q (a^\top u_i + b) - \sum_{i=1}^m \log \left(1 + \exp(a^\top u_i + b) \right).$

MLE of a, b is $\ell(a, b)$ $-\log(\cdot)$: decreasing convex.

$(\hat{a}_{ML}, \hat{b}_{ML}) = \underset{a, b}{\operatorname{argmax}} \sum_{i=1}^q (a^\top u_i + b) - \sum_{i=1}^m \log \left(1 + \exp(a^\top u_i + b) \right)$

Concave in a, b . Convex in (a, b)

Concave in (a, b)

$\ell(a, b)$ is concave in (a, b)

Hence, the problem is a convex opt. problem & can be solved numerically & efficiently.

Covariance Estimation of Gaussian random vectors.

$y \in \mathbb{R}^n$: Gaussian random vector. unknown

zero-mean $E[y] = \underline{0}$, covariance matrix $E[yy^\top] = R$.

Probability density function:

$$p_R(y) = \frac{1}{(2\pi)^{n/2} |R|^{1/2}} \exp(-y^\top R^{-1} y / 2), \quad R \in S^n_+.$$

$$|R| = \det(R)$$

Estimate R from N independent samples $y_1, \dots, y_N \in \mathbb{R}^n$.

$$\lambda_y(R) = \log p_R(y_1, \dots, y_N)$$

$$= \log \left[\prod_{k=1}^N \frac{1}{(2\pi)^{n/2} |R|^{1/2}} \exp(-y_k^\top R^{-1} y_k / 2) \right].$$

$$= -\frac{Nn}{2} \log(2\pi) - \frac{N}{2} \log \det(R) - \frac{1}{2} \sum_{k=1}^N y_k^\top R^{-1} y_k$$

Define $\bar{Y} = \frac{1}{N} \sum_{k=1}^N y_k y_k^\top$ sample covariance matrix of observations y_1, \dots, y_N . Based on this, want to estimate R

$$\lambda_y(R) = -\frac{Nn}{2} \log(2\pi) - \frac{N}{2} \log \det(R) - \frac{N}{2} \text{tr}(R^{-1} \bar{Y})$$

$$\frac{1}{2} \sum_{k=1}^N y_k^\top R^{-1} y_k = \frac{N}{2} \left(\frac{1}{N} \sum_{k=1}^N y_k^\top R^{-1} y_k \right)$$

$$= \frac{N}{2} \frac{1}{N} \sum_{k=1}^N \text{tr}(y_k^\top R^{-1} y_k) \quad |X| \text{ matrix}$$

the trace is itself

$$= \frac{N}{2} \frac{1}{N} \sum_{k=1}^N \text{tr}(R^{-1} y_k y_k^\top) = \frac{N}{2} \text{tr}(R^{-1} \frac{1}{N} \sum_{k=1}^N y_k y_k^\top)$$

$$= \frac{N}{2} \text{tr}(R^{-1} Y).$$

need this to be concave
but it is convex!

$$\max_{R \in S_{\text{ff}}^n} l_R(R) = -\frac{Nn}{2} \log(2\pi) - \frac{N}{2} \log \det(R) - \frac{N}{2} \text{tr}(R^{-1} Y)$$

$$\det(S^{-1}) = \frac{1}{\det(S)}.$$

This is not a convex opt. problem in R .

Let $S = R^{-1}$ (precision/information/inverse covariance matrix).

$$\max_{S \in S_{\text{ff}}^n} l_S(S) = -\frac{Nn}{2} \log(2\pi) + \frac{N}{2} \log \det(S) - \frac{N}{2} \text{tr}(SY).$$

(inner product linear)

$\log \det(S)$ is concave in S & $\text{tr}(SY)$ is also concave in S .

Hence the ML problem for S is a convex optimization problem.

$$S_{\text{ML}} = \underset{S \in \mathcal{A}}{\operatorname{argmax}} \log \det(S) - \text{tr}(SY).$$

convex opt prob

\mathcal{A} : encode some prior knowledge on S .

$R-L$ is +ve semi definite

i) $L \leq R \leq U \Leftrightarrow U^{-1} \leq S \leq L^{-1}$

ii) $\frac{\lambda_{\max}(R)}{\lambda_{\min}(R)} \leq K$: condition number

$$\Leftrightarrow \frac{\lambda_{\max}(S)}{\lambda_{\min}(S)} \leq K$$

iii) S to be sparse matrix \Rightarrow Gaussian graphical models

$$S_{ML} = \underset{S \in \mathcal{A}}{\operatorname{argmax}} \log \det(S) - \operatorname{tr}(SY) - \lambda \|S\|_1$$

inverse covariance matrix sparse
by adding regularization parameter

$$\|S\|_1 = \sum_{i,j} |S_{ij}| \text{ induces sparsity in matrix } S.$$

Experiment Design

Regression: Estimate unknown parameter vector $x \in \mathbb{R}^n$ given some measurements or experiments.

$$y_i = \left(\begin{matrix} a_i^T \\ \vdots \end{matrix} \right) \text{ feature vectors} \quad \text{choose from a set?} \\ y_i = a_i^T x + w_i, \quad i=1, \dots, m.$$

w_i : measurement noise $w_i \sim N(0, 1)$.

Assume a_1, \dots, a_m span \mathbb{R}^n . Set

$$A = \begin{bmatrix} -a_1^\top \\ \vdots \\ -a_m^\top \end{bmatrix} \in \mathbb{R}^{m \times n}$$

Then we know that $A^\top A \in \mathbb{R}^{m \times m}$ is invertible.

$$A^\top A = \begin{bmatrix} 1 & \cdots & 1 \\ a_1^\top & \cdots & a_m^\top \end{bmatrix} \begin{bmatrix} -a_1^\top \\ \vdots \\ -a_m^\top \end{bmatrix} = \sum_{i=1}^m a_i a_i^\top \text{ has rank } n.$$

MLE of x : $x_{\text{me}} = \underset{x}{\operatorname{arg\min}} \quad \|Ax - y\|_2^2$ minimize sum of residual

$$\hat{x} = x_{\text{me}} = (A^\top A)^{-1} A^\top y = \left(\sum_{i=1}^m a_i a_i^\top \right)^{-1} \sum_{i=1}^m a_i y_i.$$

Estimation error covariance matrix $E[ee^\top]$ where

$$e = \hat{x} - x.$$

Error vector

$$\begin{aligned} e &= \hat{x} - x = (A^\top A)^{-1} A^\top y - x = (A^\top A)^{-1} A^\top (Ax + w) - x \\ &= \cancel{(A^\top A)^{-1}(A^\top A)x} + (A^\top A)^{-1} A^\top w - \cancel{x} \\ &= (A^\top A)^{-1} A^\top w. \end{aligned}$$

Estimation error covariance matrix $E[ee^\top]$

$$= E[(A^\top A)^{-1} A^\top w ((A^\top A)^{-1} A^\top w)^\top]$$

$$= E[(A^\top A)^{-1} A^\top w w^\top A (A^\top A)^{-1}].$$

components of
variance I
/ w indep and all have

$$= (\mathbf{A}^T \mathbf{A})^{-1} \cancel{\mathbf{A}^T} \underbrace{\mathbb{E}[ww^T]}_{\mathbf{I}} \cancel{\mathbf{A}} (\mathbf{A}^T \mathbf{A})^{-1} = (\mathbf{A}^T \mathbf{A})^{-1} = \mathbf{E}.$$

\mathbf{E} characterizes the accuracy of the estimation or the informativeness of the experiments.
 ↳ the smaller the \mathbf{E} is, the better (choose \mathbf{A} huge?)

Suppose a_1, \dots, a_m can be only chosen from p possible test vectors $v_1, \dots, v_p \in \mathbb{R}^n$.

$\underbrace{\text{can only choose from these rows}}$ $p \ll m$, error covariance should be small

Choose a_1, \dots, a_m from the set $\{v_1, \dots, v_p\}$ s.t. \mathbf{E} is "small" in some precise sense.

51	9	0	70
m_1	m_2	m_3	m_4

$m = 130$.
 how to allocate
 to choose these numbers s.t.
 \mathbf{E} is as small as possible)

Let $m_j \in \mathbb{N}$ denote the # of experiments in which a_i is chosen to be v_j .

want to choose these numbers in a more principled way
 (convex opt.)

$$\mathbf{E} = \left(\sum_{i=1}^m a_i a_i^T \right)^{-1} = 51v_1 v_1^T + 9v_2 v_2^T + 0v_3 v_3^T + 70v_4 v_4^T$$

want to minimise this $\rightarrow \mathbf{E} = \left(\sum_{j=1}^p m_j v_j v_j^T \right)^{-1}$ sum over p (p is small)

Expt design: Given the menu of possible choices of the expts v_1, \dots, v_p & the total # of expts m , choose the #s of each type of expt i.e., m_1, m_2, \dots, m_p s.t. \mathbf{E} is small.

$$m_1 + m_2 + \dots + m_p = m.$$

constraint set already non-convex

0, 1, 2, 3...

$$\min_{\text{w.r.t. } S_{ff}^n} \left(\underbrace{\left(\sum_{j=1}^p m_j v_j v_j^\top \right)^{-1}}_{\text{Estimation error}} \right) \quad \text{s.t.} \quad \sum_{i=1}^p m_i = m$$

$m_i \in \mathbb{N} \cup \{0\}$

MS are integers

Estimation error

Integer program. (not convex)

fraction of times
 v_i is used

Relaxation: Let $\lambda_i = \frac{m_i}{m}$. Then $E = \frac{1}{m} \left(\sum_{j=1}^p \lambda_j v_j v_j^\top \right)^{-1}$

$$\min_{\text{w.r.t. } S_{ff}^n} \left(E = \frac{1}{m} \left(\sum_{j=1}^p \lambda_j v_j v_j^\top \right)^{-1} \right) \quad \text{matrix}$$

$$\lambda = (\lambda_1, \dots, \lambda_p) \geq 0, \quad \sum_{i=1}^p \lambda_i = 1 \Leftrightarrow 1^\top \lambda = 1$$

need to be converted into a function to be minimised

Rmk: Almost equivalent to the integer program when m is large.

Scalarization: D-optimal design ($D = \det(X)$).

$$\min_{\substack{\lambda \geq 0, 1^\top \lambda = 1 \\ \text{set of all probability vectors (convex)}}} \log \det \left(\sum_{j=1}^p \lambda_j v_j v_j^\top \right)^{-1}$$

$\log \det(\frac{1}{\lambda}) = -\log \lambda$
which is convex

This is a convex optimization problem.

Rewrite as:

$$\begin{aligned} & \min_{\lambda \geq 0, 1^\top \lambda = 1} Z \\ & \quad X = \sum_{j=1}^p \lambda_j v_j v_j^\top \end{aligned}$$

introduce X as a new parameter

Lagrangian: / Lagrange multipliers

Z (capital) matrix used as L multiplier

$$L(\lambda, X; z, \mu) = \log \det(X^{-1}) + z^T(-\lambda) + \mu(1^T \lambda - 1)$$

$$+ \text{tr} \left(Z \left(X - \sum_{j=1}^p \lambda_j v_j v_j^T \right) \right)$$

Primal variables

$$= \log \det(X^{-1}) + \text{tr}(ZX) - \sum_j \lambda_j v_j^T Z v_j - z^T \lambda + \mu(1^T \lambda - 1)$$

linear

$$= \log \det(X^{-1}) + \text{tr}(ZX) + \sum_{j=1}^p \lambda_j (-v_j^T Z v_j - z_j + \mu) - \mu.$$

$$g(z, \mu) = \inf_{\lambda, X} L(\lambda, X; z, \mu)$$

= 0

The min over λ is unbounded below unless $\mu - v_j^T Z v_j = z_j$.

Differentiate w.r.t. X and set to 0. This yields $X^{-1} = Z$.

$$h(X) = \log \det(X^{-1}) + \text{tr}(ZX)$$

$$\nabla h(X) = -X^{-1} + Z = 0 \Rightarrow X^{-1} = Z.$$

Dual function: $g(z, \mu) = \begin{cases} \log \det Z + n - \mu, & \mu - v_j^T Z v_j = z_j \\ -\infty & \text{else.} \end{cases}$

Dual problem: $\max_{z, \mu} \log \det Z + n - \mu \text{ s.t. } \mu - v_j^T Z v_j = z_j$

$$Z > 0 \text{ so } \mu \geq v_j^T Z v_j$$

$$\max_{Z, \mu} \log \det Z + n - \mu \text{ s.t. } v_j^T Z v_j \leq \mu \quad \forall j = 1, \dots, p$$

$\det(\mu W) \stackrel{!!}{=} \mu^n \det(W)$

$$v_j^T \frac{Z}{\mu} v_j \leq 1$$

Eliminate μ by making a change of variables $W = \frac{1}{\mu} Z$.

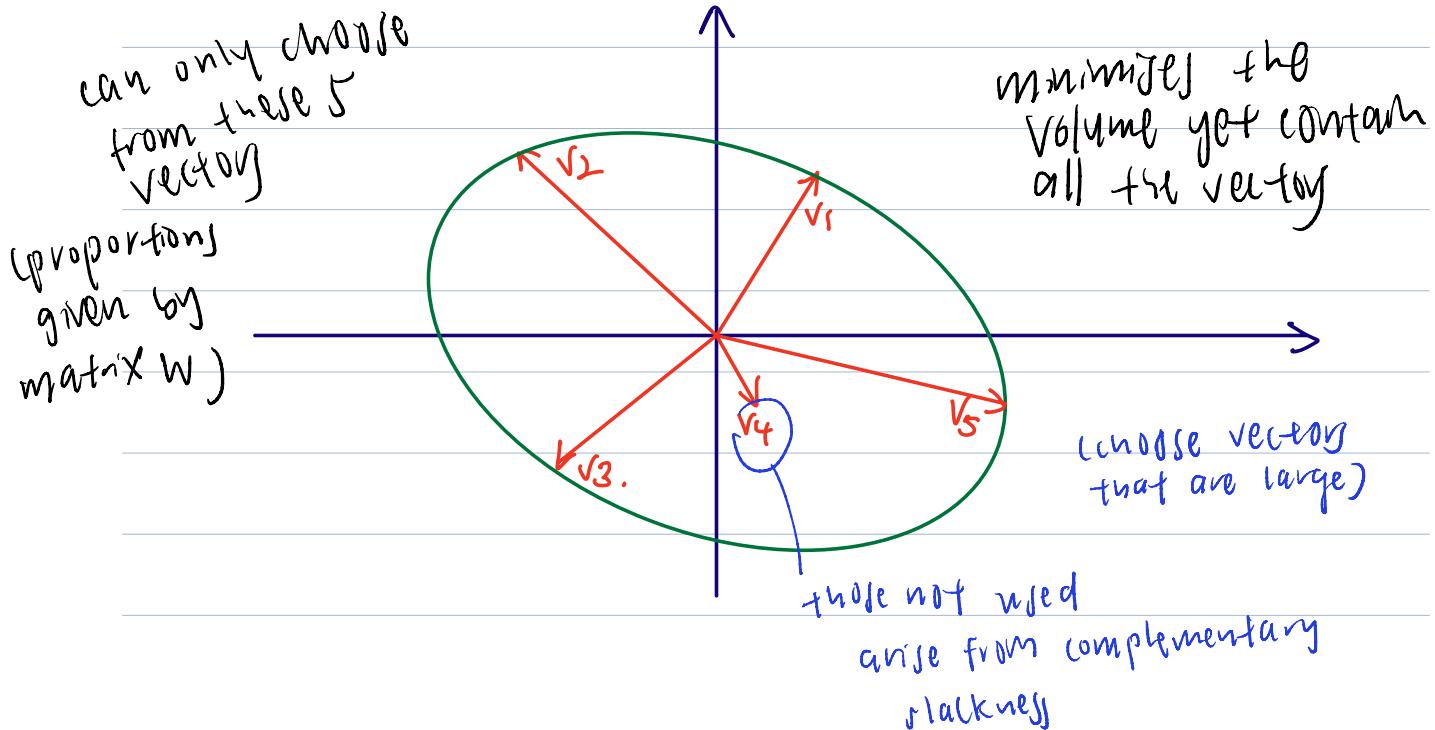
can split them up and do separately

$$\max_{W, \mu} \log \det W + n + \underline{\log \mu - \mu} \quad \text{s.t. } v_j^T W v_j \leq 1 \quad \forall j.$$

Analytically max over μ & the optimum is at $\mu = n$.

$$\max_{W \in S^n_{++}} \log \det W + n \log n \quad \text{s.t. } v_j^T W v_j \leq 1 \quad \forall j.$$

Optimal solution W^* determines the min volume ellipsoid, centered at the origin $\{x : x^T W^* x \leq 1\}$ that contains the points v_1, \dots, v_p .



dual optimal variable

By complementary slackness $(\gamma_i^*)[1 - v_i^\top W^* v_i] = 0 \quad \forall i=1,\dots,p$

i.e., the optimal design only uses expts that lie on the surface of the min volume ellipsoid.