

- We will look at designs based on a specific criteria: minimizing variance in model parameters
- Specifically, we look at couple of tools related to the previous lecture on information theory.

### Maximum entropy

- Suppose we have a model whose parameters are  $\beta$  and we specified a prior  $p(\beta)$

We can derive our posterior to be

$$p(\beta | y, D) = \frac{p(y | \beta, D) p(\beta)}{p(y | D)}$$

- Observe that we have a prior that is a simple known distribution and are interested in modeling the posterior after observing some data  $D$  that predicts distribution on outputs  $y$
- From the previous lecture, we can define a score function using the KL divergence: posterior relative to the prior

$$\text{score} = \int p(\beta | y, D) \log \left( \frac{p(\beta | y, D)}{p(\beta)} \right) d\beta dy$$

- if we are interested collecting data such that your prior and posterior look similar such that you can use one as a proxy for other, we need to minimize the score.

$$\begin{aligned} \text{score} &= \int p(\beta | y, D) \log(p(\beta | y, D)) - \underbrace{\int p(\beta | y, D) \log(p(\beta))}_{\text{constant}} \\ &= \min \{ H(p(\beta | y, D)) \} \\ &= \max \{ H(p(\beta)) \} \end{aligned}$$

(Thus fitting your posterior to the prior is maximizing the entropy of obtaining parameters  $\beta$  from posterior)

Example: one dimensional linear regression

$$\text{model: } y = A^T x + \varepsilon \sim N(A^T x, \sigma^2)$$

$$p(y | x, \beta) = N(A^T x, \sigma^2) \quad \beta = [A, \sigma^2]$$

$$\begin{aligned} \text{posterior } p(\beta | y, x) &= p(y | x, \beta) p(\beta) \\ &\propto \prod_{i=1}^n p(y_i | x_i, \beta) \\ &= \prod_{i=1}^n N(A^T x_i, \sigma^2) \\ &\propto \exp \left( -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - A^T x_i)^2 \right) \\ &= \exp(-T) \end{aligned}$$

$$\begin{aligned} \text{score} &= -\mathbb{E} [\log p(\beta | y, x)] \\ &= -\mathbb{E} \left[ \log \left( \frac{1}{(2\pi\sigma^2)^n} \exp \left( -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - A^T x_i)^2 \right) \right) \right] \\ &= -\mathbb{E} \left[ \log \left( \frac{1}{(\sqrt{2\pi\sigma^2})^n} \right) \right] - \mathbb{E} \left[ \frac{-1}{2\sigma^2} \sum_{i=1}^n (y_i - A^T x_i)^2 \right] \\ &= -\mathbb{E} \left[ \log \frac{1}{(\sqrt{2\pi\sigma^2})^n} \right] + \frac{1}{2\sigma^2} \mathbb{E} \left[ \sum_{i=1}^n (y_i - A^T x_i)^2 \right] \\ &= -\log \frac{1}{(\sqrt{2\pi\sigma^2})^n} + \frac{n}{2} \\ &= \frac{n}{2} \log(2\pi\sigma^2) + \frac{n}{2} \end{aligned}$$

### Fisher information

- consider the case of an experiment making nanoparticles of two morphologies: sphere and cylindrical
- our model for a unit volume we produce number of particles in the following way.

$$n_{\text{sph}} = \alpha + \beta \quad \text{and } \sigma_s \text{ variance}$$

$$n_{\text{cyl}} = \alpha \quad \sigma_c \text{ variance}$$

- our goal is to predict how well an experimental design will be able to constrain the model parameters before doing the experiment

- such that we can forecast the results of different volumes to measure and compare precision versus cost

- let us try to predict the output of our experiments based on some intuition.

given data, we can infer  $n_{\text{cyl}}$ ,  $n_{\text{sph}}$  by using the model above.

if there are too many  $n_{\text{cyl}}$  in the volume that we measured

estimate of  $\alpha$  but depress the estimate  $\beta$

- That means there is a covariance between our estimate of the two parameters

- variance in  $\sigma_s$  is the sum of variance in  $\alpha$  and  $\beta$  while  $\sigma_c$  is equal to variance in  $\alpha$ .

- The Fisher information matrix defined below makes this idea precise and quantitative

we define our mode to be of  $N$  parameters

$p_1, p_2, \dots, p_N$  such that our observables are

given by

$$f_b = f_b(p_1, p_2, \dots, p_N)$$

$$\text{Fisher information matrix } F_{ij} = \sum_b \frac{1}{\sigma_b^2} \frac{\partial f_b}{\partial p_i} \frac{\partial f_b}{\partial p_j}$$

- For the above simple model, we get

$$F = \begin{bmatrix} \frac{\alpha}{\sigma_s^2} + \frac{1}{\sigma_c^2} & \beta \\ \beta & \frac{1}{\sigma_c^2} \end{bmatrix}$$

covariance of the model parameters is  $F^{-1}$

$$= \begin{bmatrix} \sigma_s^2 & -\sigma_s^2 \\ -\sigma_s^2 & \sigma_s^2 + \sigma_c^2 \end{bmatrix}$$

$$p_1 = \alpha, \quad p_2 = \beta$$

$$n_s = f_1, \quad n_{\text{cyl}} = f_2$$

$$F_{1,1} = \frac{1}{\sigma_s^2} \frac{\partial f_1}{\partial \alpha} \frac{\partial f_1}{\partial \alpha} + \frac{1}{\sigma_c^2} \frac{\partial f_2}{\partial \alpha} \frac{\partial f_2}{\partial \alpha}$$

$$= \frac{1}{\sigma_s^2} \left( 1 + \frac{\partial f_1}{\partial \alpha} \right)^2 + \frac{1}{\sigma_c^2} = \frac{1}{\sigma_s^2} + \frac{1}{\sigma_c^2}$$

- which verifies that  $\alpha$  (# of cylinders in vol.) is proportional to variance in cylinder production,

- one of the advantage of using Fisher information matrices is that we can encode prior knowledge by adding appropriate variance term to the matrix

for example if we know that # of spherical particles in a given volum has the variance  $\sigma_{s,p}^2$  then the fisher matrix should be modified as

$$\text{prior} = \begin{bmatrix} \sigma_{s,p}^2 & 0 \\ 0 & \sigma_{c,p}^2 \end{bmatrix}$$

$$I = \text{prior}^{-1}$$

$$\text{thus } F_{\text{prior}} = F + \text{prior}^{-1}$$

(show example from dexp)

- one algorithmic way of achieving optimal design is to use the coordinate-exchange algorithm.

- at each iteration we randomly select two points from the candidate set, swap their values and see if the criteria improves. Repeat until we don't see any improvement

- The exchange step is primarily to improve exploration and avoid local optimum

- Another simple example would be measuring tensile strength by application of different levels of stress. Our regression model is more precise when the determinant of Fisher information matrix is large (low error of the model)

(For any matrix, its determinant computes the volume defined by cuboid it spans when applied to any vector of the corresponding)

using this, we would say that D-optimality tries to maximize volume of information

(show example from dexp)

- one algorithmic way of achieving optimal design is to use the coordinate-exchange algorithm.

- at each iteration we randomly select two points from the candidate set, swap their values and see if the criteria improves. Repeat until we don't see any improvement

- The exchange step is primarily to improve exploration and avoid local optimum

- Another simple example would be measuring tensile strength by application of different levels of stress. Our regression model is more precise when the determinant of Fisher information matrix is large (low error of the model)

(For any matrix, its determinant computes the volume defined by cuboid it spans when applied to any vector of the corresponding)

using this, we would say that D-optimality tries to maximize volume of information

(show example from dexp)