

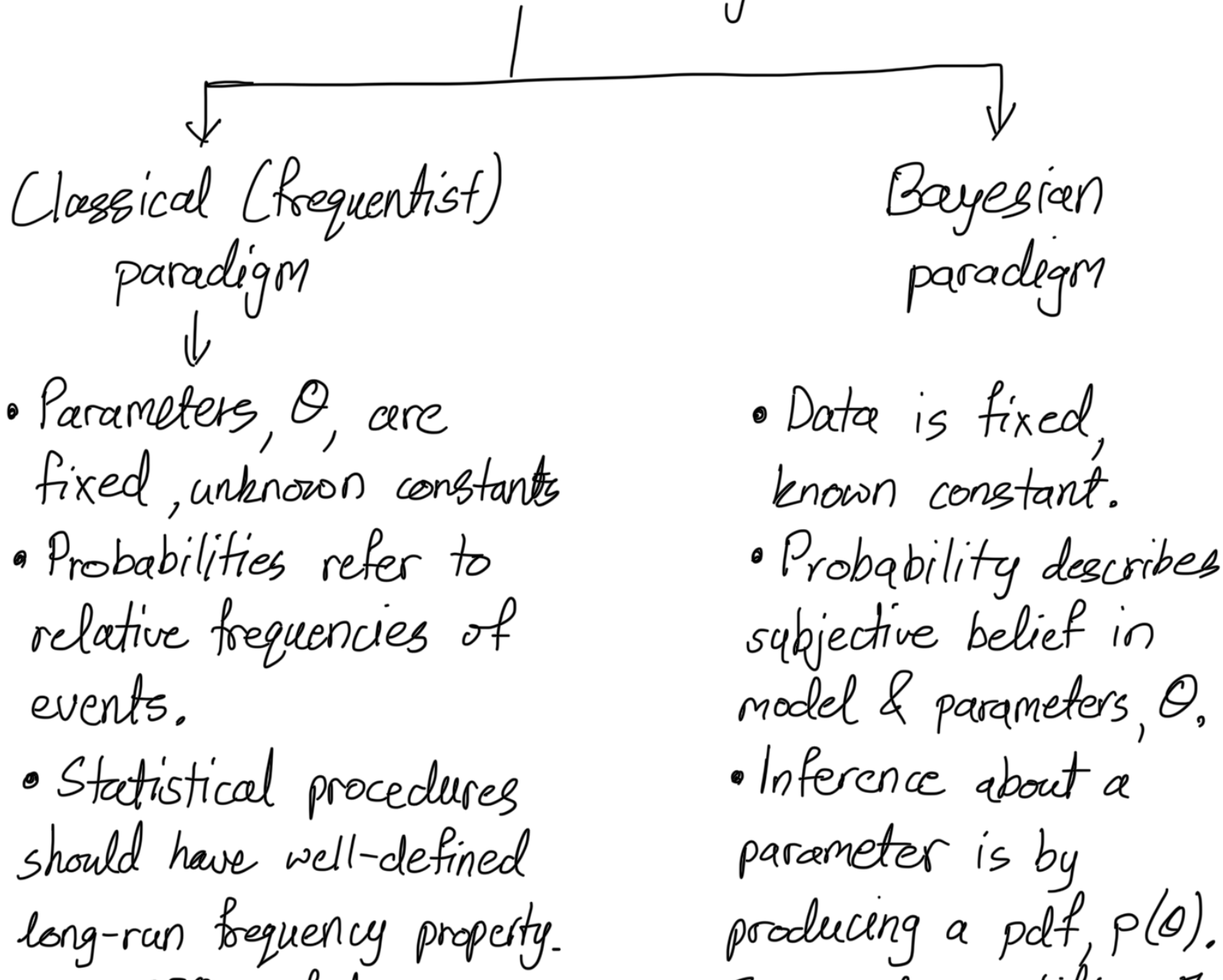
Lecture 5

Classical Statistical Inference

* Three main concepts in statistical inference:

- 1) Point estimate: what is the best estimate for a model parameter θ .
- 2) Confidence estimate: How confident should we be in that estimate?
- 3) Hypothesis testing: Are the data consistent with the model used?

Statistical Paradigms



eg: a 95% confidence interval should bracket the true value with a limiting frequency of at least 95%.

- Can be computationally cheap/easy.

This plot quantifies the uncertainty about the parameter & can be used to make point estimate

- Usually more computationally demanding.

- The important thing to note that the best-fit solutions are typically consistent between frequentist & Bayesian methods.

Maximum Likelihood Estimation

- MLE is a common special case of frequentist analyses.

- However, Bayesian analyses build on top of the MLE framework too.

- Note: frequentist analyses, unlike Bayesian analyses, are not tied to the likelihood.

- You can use other cost functions in frequentist analyses, eg. MSE, ELU, RELU.

- The data "likelihood" represents a quantitative description of our measuring process.

eg. if we know (or assume) that our data $\{x_i\}$ are drawn from $N(\mu, \sigma)$, then the likelihood for any single data point is

$$p(x_i | \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

$$p(x_i | \mu, \sigma) = N(x_i | \mu, \sigma)$$

- For the full dataset, the likelihood is given by a product of individual probabilities / lkl:

$$\mathcal{L} \equiv p(\{x_i\} | M(\vec{\theta})) = \prod_{i=1}^N p(x_i | M(\vec{\theta}))$$

- If our data is drawn from a Gaussian, then $M(\vec{\theta}) \equiv [\mu, \sigma]$.

- Note that the likelihood, \mathcal{L} , is not a true (properly normalized) pdf, nor is it a statement on the probability of obtaining model parameters, $\vec{\theta}$.

- Since \mathcal{L} is a product of probabilities, all < 1 , the \mathcal{L} can be a very small number, and we instead work with $\ln \mathcal{L}$ instead.

- The maximum likelihood approach can be broken down into the following steps:

① Formulate the data likelihood, $p(D|M)$, for some model $M(\vec{\theta})$ with parameters $\vec{\theta}$.

② Search for the parameters, $\vec{\theta}^*$, that maximize $p(D|M)$.

- These parameters, $\vec{\theta}^*$, are the point estimates.

③ Determine the confidence region for model parameters $\vec{\theta}^*$:

- either analytically with the Fisher matrix if \mathcal{L} is well-behaved.

- numerically by either bootstrap, jack-knife or cross-validation.

④ Perform hypothesis testing to see if your chosen model explains the observed data. If not, find a new $M'(\vec{\theta})$ and go back to ①.

- Properties of ML estimators:

- They are consistent estimators, i.e. they can be proven to converge to the true parameter value as number of data points increases.
- They are asymptotically normal estimators, i.e. as number of data points increase, the parameter distribution approaches a normal distribution centered on the MLE.
- They asymptotically achieve the theoretical minimum possible variance, called the Cramer-Rao bound.

- Thus, to summarize, given a likelihood, L , the MLE, $\vec{\theta}^0$, is given by:

$$\left. \frac{d \ln L(\vec{\theta})}{d \vec{\theta}} \right|_{\vec{\theta}^0} = 0$$

- Using the asymptotic normality of MLE, the uncertainty matrix for these parameters can be derived by,

$$\sigma_{jk} = \left([F^{-1}]_{jk} \right)^{1/2}$$

where,

$$F_{jk} = - \frac{d^2 \ln \mathcal{L}}{d\theta_j d\theta_k} \bigg|_{\theta = \theta_0}$$