# Lecture #5

· The central idea of statistics is to infer from data.

· That is, let $X$ be some unknown r.v. If we observe $X_1, X_2, ..., X_n \overset{iid}{\sim} X$, what can we say about $X$ itself?

· Possible questions about $X$:
- What is the density of $X$? This is often very hard
- Does $\underline{X}$ look like it comes from a certain family of distributions?
- What is $\mathbb{E}(X)$? $Var(X)$? More generally,
$$\int_{\mathbb{R}} g(x) \cdot f_X(x) dx \quad \text{for some } g(x).$$

· Two broad directions: parametric v. nonparametric

            ↓                     ↓

    ~ makes strong assumptions     ~ few assumptions → very flexible

    ~ easy to work with         ~ more difficult computationally

<u>Def'n</u>: A collection of functions $\mathcal{F}$ is a parametric model if it can be parametrized by a finite number of parameters i.e. ~~xxxxx~~ There exists a finite dimensional parameter space $\Theta = \{(\Theta_1, \Theta_2, ..., \Theta_D)\}$ s.t. every element of $\mathcal{F}$ is associated with an element of $\Theta$.

ex: Everything we ~~xxx~~ saw in MATH 165 (more or less).

ex: Normal distributions in $\mathbb{R}$ is a parametric model with two parameters:

$f(x)$ is a normal density if and only if there exist $\mu \in \mathbb{R}$ ("mean")

$\sigma^2 \in \mathbb{R}_+$ ("variance")

s.t. $\quad f(x) = \cancel{\text{\#\#\#}} \dfrac{1}{\sqrt{2\sigma^2}} \cdot \exp\left(\dfrac{-[x-\mu]^2}{2\sigma^2}\right).$

· Parametric models make identifying a good estimate for the true density easier... just estimate the finitely many parameters! How hard can that be...

· If we consider $\mathcal{F}$ to not have a finite parametrization, this is non-parametric statistics.

ex: Let $\mathcal{F} = \left\{ f \in \mathcal{C}(\mathbb{R}) \,\middle|\, f(x) \geq 0 \text{ and } \int_{\mathbb{R}} f(x)\,dx \overset{=}{\phantom{1}} 1 \right\}$ be the space of all continuous densities. This is a big space of functions, and doesn't have a useful finite parametrization.

· Context: Given data $\{x_i\}_{i=1}^n \overset{iid}{\sim} \underset{\sim}{X}$ for $\underset{\sim}{X}$ unknown,

  (1) Parametric: Find the best Gaussian to model $\underset{\sim}{X}$

  (2.) Nonparametric: Find the best continuous density to model $\underset{\sim}{X}$

· Obviously (2.) is less constrained.

- An underlying issue with statistics is how to decide if I have predicted well.

- That is, given data $X_1, ..., X_n \overset{iid}{\sim} X$, how do I know if my predicted density or predicted ~~XXXX~~ $\mathbb{E}(X)$ is good or not? This is a problem that never really goes away, and can be studied on many levels.

- We will focus on the "classical" approach to understanding error in statistical estimation: bias-variance tradeoff.

/

- Density estimation is hard, and we will focus on an easier problem for now: <u>point estimation</u>, i.e. estimating a particular quantity of interest. We are especially interested in estimating the expected value of $X$, given only observations $X_1, ..., X_n \overset{iid}{\sim} X$.

- More generally, let $\Theta$ be an underlying quantity of interest, quite of a parameter in a parametric model (e.g. $\Theta = \mathbb{E}(X)$). $\Theta$ itself is unknowable, so we would like to estimate it given (random) data $X_1, ..., X_n \sim X$.

- We can think of a method of estimation as a <u>function on the data</u>. This is a fundamental insight of statistics, because it allows for the machinery of analysis and probability to give us control.

- So, we will write $\hat{\Theta}(X_1, ..., X_n)$ as a function that estimates $\Theta$, given

data $X_1, \ldots, X_n \sim^{iid} X$.

ex: Let $X_1, \ldots, X_n \sim X$, where $\mathbb{E}(X)$ exists. A classical problem is to estimate $\mathbb{E}(X)$. It is almost an article of faith that we should use $\frac{1}{n} \sum_{i=1}^{n} \hat{X}_i$.

In 1-dimension, a quick simulation shows that to be effective for many $X$, if $n$ is large enough. 

In fact the WLLN says this is a good idea! ~~[scribbled out]~~

Indeed, WLLN says, $\forall \varepsilon > 0$,

$$\lim_{n \to \infty} \mathbb{P}\left( \left| \frac{1}{n} \sum_{i=1}^{n} X_i - \mathbb{E}(X) \right| > \varepsilon \right) = 0.$$

So, if $\left. \begin{array}{l} \theta = \mathbb{E}(X) \\ \hat{\theta}(x_1, \ldots, x_n) = \frac{1}{n} \sum_{i=1}^{n} X_i, \\ \quad\quad \hat{\theta}_n \end{array} \right\} \xrightarrow{WLLN} \begin{array}{l} \forall \varepsilon > 0, \\ \lim_{n \to \infty} \mathbb{P}\left( |\hat{\theta}_n - \theta| > \varepsilon \right) = 0. \\ \quad\quad\quad \text{pretty good!} \end{array}$

This motivates the idea of consistency.

Defn: ~~[scribbled]~~ A point estimator $\hat{\theta}(x_1, \ldots, x_n) = \hat{\theta}_n$ is <u>consistent</u> if

$$\hat{\theta}_n \xrightarrow{P} \theta.$$

- A key insight is that $\hat{\theta}_n$ is a r.v. So, we can consider its expectation and variance.

**Def^n:** Let $\hat{\Theta}_n$ be an estimator of $\Theta$. ~~[scribbled out]~~ The ~~bias~~ bias of $\hat{\Theta}_n$ is $\text{bias}(\hat{\Theta}_n) = \mathbb{E}(\hat{\Theta}_n) - \Theta$, where the expectation is taken over $X_1, \ldots, X_n \sim X$. We say $\hat{\Theta}_n$ is <u>unbiased</u> if $\text{bias}(\hat{\Theta}_n) = 0$.

· Let $\text{Var}(\hat{\Theta}_n) = \mathbb{E}\left(\left[\hat{\Theta}_n - \mathbb{E}(\hat{\Theta}_n)\right]^2\right)$ be the variance of $\hat{\Theta}_n$ with respect to the sample $X_1, \ldots, X_n \sim X$. We are typically interested in understanding the "total" error of estimating:

$$MSE(\hat{\Theta}_n) = \mathbb{E}\left(\left[\hat{\Theta}_n - \Theta\right]^2\right),$$ the mean square error.

· Amazingly, this is interpretable into two components = model part v. random part.

<u>Theorem</u> (Bias-Variance Tradeoff) : $MSE(\hat{\Theta}_n) = \text{bias}(\hat{\Theta}_n)^2 + \text{Var}(\hat{\Theta}_n)$.

**Proof:** $\underbrace{\mathbb{E}\left(\left[\hat{\Theta}_n - \Theta\right]^2\right)}_{}$ $\quad {}_{=A} \quad {}_{=B} \quad$ $\mathbb{E}\left(\left[A+B\right]^2\right) = \mathbb{E}A^2 + \mathbb{E}2AB + \mathbb{E}B^2$

$= \mathbb{E}\left(\left[\overbrace{\hat{\Theta}_n - \mathbb{E}(\hat{\Theta}_n)}^{A} + \overbrace{\mathbb{E}(\hat{\Theta}_n) - \Theta}^{B}\right]^2\right)$

$= \underbrace{\mathbb{E}\left(\left[\hat{\Theta}_n - \mathbb{E}(\hat{\Theta}_n)\right]^2\right)}_{\overset{=}{Var(\hat{\Theta}_n)}} + 2\,\mathbb{E}\left(\left[\hat{\Theta}_n - \mathbb{E}(\hat{\Theta}_n)\right]\overbrace{\left[\mathbb{E}(\hat{\Theta}_n) - \Theta\right]}^{\text{constant}}\right) + \mathbb{E}\left(\overbrace{\left[\mathbb{E}(\hat{\Theta}_n) - \Theta\right]^2}^{\text{constant}}\right)$

$= Var(\hat{\Theta}_n) + 2\left(\mathbb{E}(\hat{\Theta}_n) - \Theta\right) \cdot \underbrace{\mathbb{E}\left[\hat{\Theta}_n - \mathbb{E}(\hat{\Theta}_n)\right]}_{= 0} + \underbrace{\left[\mathbb{E}(\hat{\Theta}_n) - \Theta\right]^2}_{= \text{bias}(\hat{\Theta}_n)^2}$

$= Var(\hat{\Theta}_n) + \text{bias}(\hat{\Theta}_n)^2.$ ∎

· $\text{bias}(\hat{\Theta}_n) \approx 0$ if there is enough capacity in the model to fit it well. $Var(\hat{\Theta}_n) \approx 0$ if we have enough data to reliably estimate.