

Lecture # 7

Let us now consider parametric models. We consider a space of functions parametrized by $\Theta \subset \mathbb{R}^D$: $\mathcal{F} = \{f(x; \theta) \mid \theta \in \Theta\}$. We can think of $f(x; \theta)$ two ways:

(1.) For fixed θ , $x \mapsto f(x; \theta)$ can be the density/cdf associated to a fixed parameter set.

(2.) For fixed x , $\theta \mapsto f(x; \theta)$ illustrates how the family changes with respect to the parameters.

A central question for us is, given data, how to determine a "good" choice of parameter? We'll focus on two high level methods: (1.) Method of moments (MM) (2.) Maximum likelihood estimation (MLE)

~~Let us~~ Let us recall the definition of moments of a r.v.

Defn. Let X be a r.v. (1.) ~~For~~ For $p \geq 1$, the p^{th} moment of X is $E(X^p)$.

(2.) For $p \geq 1$, the p^{th} sample moment of X is

$$\frac{1}{n} \sum_{i=1}^n X_i^p, \quad \text{where } \{X_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} X.$$

The idea of MM is to choose parameters so that the exact moment matches the sample moment.

②
That is, let $X(\theta)$ be the r.v. associated to parameter θ . Then we can compute the (exact) p^{th} moment $\alpha_p(\theta) = E(X(\theta)^p)$.

Given data, we can compute sample moments $\hat{\alpha}_p = \frac{1}{n} \sum_{i=1}^n X_i^p$, where $\{X_i\}_{i=1}^n \stackrel{\text{iid}}{\sim} X$.

The MM idea is to put these in correspondence via a system of equations: choose $\hat{\theta} \in \Theta \subseteq \mathbb{R}^D$ s.t. $\alpha_p(\hat{\theta}) = \hat{\alpha}_p$, $p=1, \dots, D$.

Remark - We need to do this for $p=1, \dots, D$ in order to have D equations in the D unknowns of Θ . This can get a little tedious if D is large!

ex: Consider a 1 parameter family of exponentials: $f_\lambda(x) = \begin{cases} \lambda e^{-\lambda x} & , x \geq 0 \\ 0 & , \text{else} \end{cases}$

Given data $\{X_i\}_{i=1}^n \stackrel{\text{iid}}{\sim} X$, our task is to estimate a good choice of λ . Since there is only one parameter (λ), MM says "choose λ^* so that f_{λ^*} has E matching the empirical average observed from the data."

So, choose λ^* s.t. $E(\bar{X}(\lambda)) = \int_0^\infty x \cdot \lambda e^{-\lambda x} dx = \frac{1}{n} \sum_{i=1}^n x_i$.

Through IBP, we see $\int_0^\infty x \cdot \lambda \exp(-\lambda x) dx = \frac{1}{\lambda}$, so we simply set

$\lambda^* = \left(\frac{1}{n} \sum_{i=1}^n x_i \right)^{-1}$ Easy enough!

ex: Let us consider the 2-parameter family of normal distributions, i.e.

$\left\{ \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{[x-\mu]^2}{2\sigma^2}\right) \right\}_{\substack{\mu \in \mathbb{R}, \\ \sigma \in \mathbb{R}_+}}$ Given data $\{x_i\}_{i=1}^n \sim \mathcal{N}$ (unknown \bar{X}), our

goal is to estimate "good" parameters σ^2 and μ . The MM suggests choosing

$(\hat{\mu}, \hat{\sigma}^2)$ s.t. $E(\bar{X}(\hat{\mu}, \hat{\sigma}^2)) = \frac{1}{n} \sum_{i=1}^n x_i$

$E(\bar{X}(\hat{\mu}, \hat{\sigma}^2)^2) = \frac{1}{n} \sum_{i=1}^n x_i^2$, $\mathcal{N}(\mu, \sigma^2)$ the Gaussian r.v. with parameters (μ, σ^2) .

We know $E(\bar{X}(\hat{\mu}, \hat{\sigma}^2)) = \hat{\mu}$, so we know the MM choice is

$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$. Natural enough!

For the second moment, we recall $\text{Var}(\bar{X}) = E(\bar{X}^2) - E(\bar{X})^2$ for any \bar{X} ,

and hence $E(X(\hat{\mu}, \hat{\sigma}^2)^2) = \text{Var}(X(\hat{\mu}, \hat{\sigma}^2)) + E(X(\hat{\mu}, \hat{\sigma}^2))^2$ (4)

$$= \hat{\sigma}^2 + \hat{\mu}^2.$$

So, MM requires us to solve the system
$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\hat{\mu} + \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2$$

$$\Rightarrow \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \frac{1}{n} \sum_{i=1}^n x_i.$$

Again, natural enough.

The MM has some theoretical guarantees we won't prove:

Theorem (Properties of MM): Let $\hat{\theta}_n$ denote the MM estimator. Under appropriate

conditions on the family $\mathcal{F} = \{f(x; \theta) \mid \theta \in \Theta\}$,

(1) $\hat{\theta}_n$ exists with probability 1 (almost surely).

(2) $\hat{\theta}_n \xrightarrow{P} \theta$.

(3) $\sqrt{n}(\hat{\theta}_n - \theta) \rightarrow \mathcal{N}(0, \Sigma)$, $\Sigma = g E(Y Y^T) g^T$, $Y = (X, X^2, \dots, X^d)$
 $g = \left(\frac{\partial \alpha_1(\theta)}{\partial \theta}, \dots, \frac{\partial \alpha_d(\theta)}{\partial \theta} \right)$

• Suffice it to say, the MM is almost surely consistent and asymptotically normal.