

# Class Notes (experimental)

*Jonathan Rosenblatt*

*April 7, 2015*

## Contents

<b>Estimation</b>	<b>1</b>
Moment matching . . . . .	2
Example: Exponential Rate . . . . .	2
Example: Linear Regression . . . . .	2
Quantile matching . . . . .	2
Example: Exponential rate: . . . . .	2
Maximum Likelihood . . . . .	2
Example: Exponential rate: . . . . .	3
Example: Discrete time Markov Chain: . . . . .	3
Example: Brownian motion with drift . . . . .	3
Example: Linear Regression . . . . .	4
M-Estimation (Empirical Risk Minimization) . . . . .	4
<b>From Estimation to Learning</b>	<b>4</b>
Empirical Risk Minimization (ERM) and Inductive Bias . . . . .	4
Linear Regression (OLS) . . . . .	4
Ridge Regression . . . . .	4
Logistic Regression . . . . .	4
LASSO . . . . .	4
Linear SVM . . . . .	4
Generalized Additive Models (GAMs) . . . . .	4
Neural Nets (NNETs) . . . . .	4
Classification and Regression Trees (CARTs) . . . . .	4

## Estimation

In this section, we present several estimation principles. Their properties are not discussed, as the section is merely a reminder and a preparation for the **Learning**.

## Moment matching

The fundamental idea: match empirical moments to theoretical. I.e., estimate

$$E[g(X)]$$

by

$$\mathbb{E}[g(X)]$$

where  $\mathbb{E}[g(X)] := \frac{1}{n} \sum_i g(X_i)$ , is the empirical mean.

### Example: Exponential Rate

Estimate  $\lambda$  in  $X_i \sim \exp(\lambda)$ ,  $i = 1, \dots, n$ , i.i.d.  $E[X] = 1/\lambda \Rightarrow \hat{\lambda} = 1/\mathbb{E}[X]$

### Example: Linear Regression

Estimate  $\beta$  in  $Y \sim \mathcal{N}(X\beta, \sigma^2 I)$ , a  $p$  dimensional random vector.  $E[Y] = X\beta$  and  $\mathbb{E}[Y] = y$ . Clearly, moment matching won't work because no  $\beta$  satisfies  $X\beta = Y$ . A technical workaround: Since  $\beta$  is  $p$  dimensional, I need to find some  $g(Y) : \mathbb{R}^n \mapsto \mathbb{R}^p$ . Well,  $g(Y) := XY$  is such a mapping. I will use it, even though my technical justification is currently unsatisfactory. We thus have:  $E[X'Y] = X'X\beta$  which I match to  $\mathbb{E}[X'Y] = X'y$ :

$$X'X\beta = X'y \Rightarrow \hat{\beta} = (X'X)^{-1}X'y.$$

## Quantile matching

The fundamental idea: match empirical quantiles to theoretical. Denoting by  $F_X(t)$  the CDF of  $X$ , then  $F_X^{-1}(\alpha)$  is the  $\alpha$  quantile of  $X$ . Also denoting by  $\mathbb{F}_X(t)$  the Empirical CDF of  $X_1, \dots, X_n$ , then  $\mathbb{F}_X^{-1}(\alpha)$  is the  $\alpha$  quantile of  $X_1, \dots, X_n$ . The quantile matching method thus implies estimating

$$F_X^{-1}(\alpha)$$

by

$$\mathbb{F}_X^{-1}(\alpha).$$

### Example: Exponential rate:

Estimate  $\lambda$  in  $X_i \sim \exp(\lambda)$ ,  $i = 1, \dots, n$ , i.i.d.

$$F_X(t) = 1 - \exp(-\lambda t) = \alpha \Rightarrow F_X^{-1}(\alpha) = \frac{-\log(1-\alpha)}{\lambda} \Rightarrow F_X^{-1}(0.5) = \frac{-\log(0.5)}{\lambda} \Rightarrow \hat{\lambda} = \frac{-\log(0.5)}{\mathbb{F}_X^{-1}(0.5)}$$

.

## Maximum Likelihood

The fundamental idea is that if the data generating process (i.e., the **sampling distribution**) can be assumed, then the observations are probably some high probability instance of this process, and not a low probability event: Let  $X_1, \dots, X_n \sim P_\theta$ , with density (or probability)  $p_\theta(X_1, \dots, X_n)$ . Denote the likelihood, as a function of  $\theta$ :  $L(\theta) : p_\theta(X_1, \dots, X_n)$ . Then  $\hat{\theta}_{ML} := \operatorname{argmax}_\theta \{L(\theta)\}$ .

**Example: Exponential rate:**

Estimate  $\lambda$  in  $X_i \sim \exp(\lambda)$ ,  $i = 1, \dots, n$ , i.i.d. Using the exponential PDF and the i.i.d. assumption

$$L(\lambda) = \lambda^n \exp(-\lambda \sum_i X_i)$$

. Using a monotone mapping such as the log, does not change the *argmax*. Denoting  $l(\theta) := \log(L(\theta))$ , we have

$$l(\lambda) = n \log(\lambda) - \lambda \sum_i X_i$$

. By differentiating and equating 0, we get  $\hat{\lambda}_{ML} = 1/\mathbb{E}[X]$ .

**Example: Discrete time Markov Chain:**

Estimate the transition probabilities,  $p_1$  and  $p_2$  in a two state,  $\{0, 1\}$ , discrete time, Markov chain where:  $P(X_{t+1} = 1|X_t = 0) = p_1$  and  $P(X_{t+1} = 1|X_t = 1) = p_2$ . The likelihood:

$$L(p_1, p_2) = P(X_1, \dots, X_n; p_1, p_2) = \prod_{t=0}^T P(X_{t+1} = x_{t+1} | X_t = x_t).$$

We denote  $n_{ij}$  the number of observed transitions from  $i$  to  $j$  and get that  $\hat{p}_1 = \frac{n_{01}}{n_{01} + n_{00}}$ , and that  $\hat{p}_2 = \frac{n_{11}}{n_{11} + n_{10}}$ .

**Remark:** Well, this is a rather artificial example, as because of the Markov property, and the stationarity of the process, we only need to look at transition events, themselves Bernoulli distributed. This example does show, however, the power of the ML method to deal with non i.i.d. samples. As does the next example.

**Example: Brownian motion with drift**

Estimate the drift parameter  $a$ , in a discrete time Gaussian process where:  $X_{t+1} = X_t + \varepsilon; \varepsilon \sim \mathcal{N}(0, \sigma^2) \Rightarrow X_{t+1}|X_t \sim \mathcal{N}(aX_t, \sigma^2)$ .

We start with the conditional density at time  $t + 1$ :

$$p_{X_{t+1}|X_t=x_t}(x_{t+1}) = (2\pi\sigma^2)^{-1/2} \exp\left(-\frac{1}{2\sigma^2}(x_{t+1} - ax_t)^2\right)$$

. Moving to the likelihood:

$$L(a) = (2\pi\sigma^2)^{-T/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{t=1}^T (x_{t+1} - ax_t)^2\right).$$

Differentiating with respect to  $a$  and equating 0 we get  $\hat{a}_{ML} = \frac{\sum x_{t+1}x_t}{\sum x_t^2}$ .

We again see the power of the ML device. Could we have arrived to this estimator by intuition alone? Hmmmm... maybe. See that  $Cov[X_{t+1}, X_t] = a Var[X_t] \Rightarrow a = \frac{Cov[X_{t+1}, X_t]}{Var[X_t]}$ . So  $a$  can also be derived using the moment matching method which is probably more intuitive.

### **Example: Linear Regression**

Estimate  $\beta$  in  $Y \sim \mathcal{N}(X\beta, \sigma^2 I)$ , a  $p$  dimensional random vector. Recalling the multivariate Gaussian PDF:

$$p_{\mu, \Sigma}(y) = (2\pi)^{-n/2} |\Sigma|^{-1/2} \exp \left( -\frac{1}{2} (y - \mu)' \Sigma^{-1} (y - \mu) \right)$$

So in the regression setup:

$$L(\beta) = p_{\beta, \sigma^2}(y) = (2\pi)^{-n/2} |\sigma^2 I|^{-1/2} \exp \left( -\frac{1}{2\sigma^2} \|y - X\beta\|^2 \right)$$

## **M-Estimation (Empirical Risk Minimization)**

## **From Estimation to Learning**

### **Empirical Risk Minimization (ERM) and Inductive Bias**

**Linear Regression (OLS)**

**Ridge Regression**

**Logistic Regression**

**LASSO**

**Linear SVM**

**Generalized Additive Models (GAMs)**

**Neural Nets (NNETs)**

**Classification and Regression Trees (CARTs)**