

2/9/15

①

STATISTICAL LEARNING/ESTIMATION (MLE, MAP, Bayesian)

① Maximum Likelihood Estimation

→ Real World Phenomenon / Sample Space

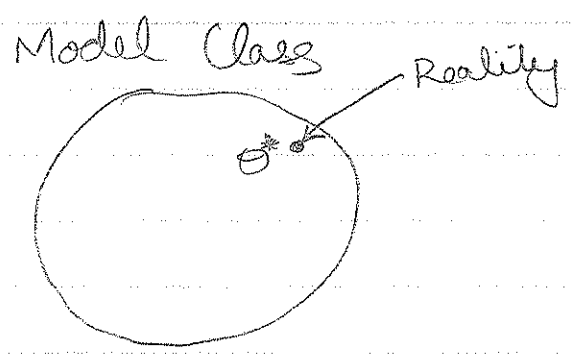
$$\Omega = \{ \text{Nadal Loses (L)}, \text{Nadal Wins (W)} \}$$

→ Random Variable: $Y = y \in \{0, 1\}$
L, W
H, T

→ Hypothesis Class / Model Class or simply "Model"

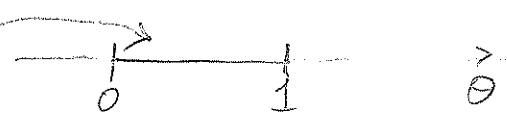
$$Y \sim \text{Bernoulli}(\theta) \iff P(Y=1) = \theta$$
$$P(Y=0) = 1 - \theta$$

Note: In this example (Nadal W/L), the model is perfect, ie no "approximation/modeling error"



[Usually, 'reality' lies outside Model Class]

As in this
Model Class



→ Learning/Estimation Goal:

Given $D = \{1, 0, 0, 1, 1\}$
estimate $\hat{\theta} \in [0, 1]$

→ How? What's a "good" θ ?

≡ best "explains" the data

≡ makes it likely for us to have observed D

[e.g. if $\theta = 0$, we would never observe a 1]

Formally, let's maximize the prob of D under θ

→ MLE

$$\hat{\theta}_{MLE} = \underset{\theta \in [0, 1]}{\operatorname{argmax}} \underbrace{P(D | \theta)}$$

called likelihood function

$$L_D(\theta) = P(D | \theta)$$

[or sometimes]
 $L(\theta; D)$]

IMP: likelihood $L_D(\theta)$ or simply $L(\theta)$ is a function of θ .

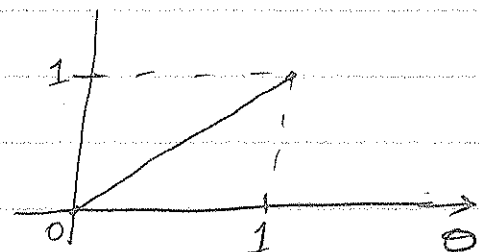
D is fixed to what we observed.

2

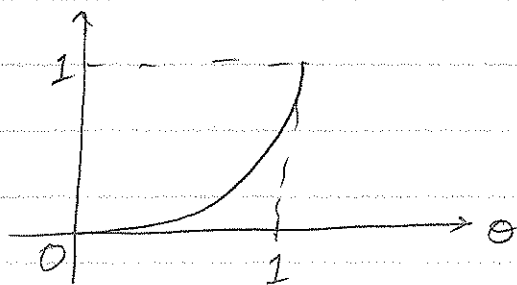
In our running example,

$$L(\theta) = P(D|\theta) = \prod_{i=1}^n P(y_i|\theta) \quad [\text{why? Hint: IID}]$$

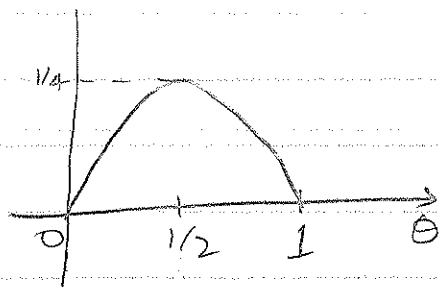
→ e.g. $D = \{1\}$ $L(\theta) = P(Y=1|\theta) = \theta$



→ $D = \{1, 1\}$ ⇒ $L(\theta) = \theta \cdot \theta = \theta^2$



→ $D = \{1, 0\}$ ⇒ $L(\theta) = P(Y=1|\theta) \cdot P(Y=0|\theta)$
 $= \theta \cdot (1-\theta)$



→ In general, $L(\theta) = \theta^{\alpha_H} (1-\theta)^{\alpha_T}$

$\alpha_H = \# \text{ Heads / Wins}$
 $\alpha_T = \# \text{ Tails / Losses}$

$$\hat{\theta}_{MLE} = \underset{\theta \in [0, 1]}{\text{argmax}} \quad L(\theta)$$

$$= \underset{\theta \in [0, 1]}{\text{argmax}} \quad \log L(\theta)$$

log-likelihood or $\log L(\theta)$

[Why? ∵ log is a monotone function, so preserves argmax.]

How do we find argmax of $LL(\theta)$?

Elementary, my dear Calculus!

Take 1st deriv; set to zero

$$\frac{\partial LL(\theta)}{\partial \theta} = \frac{\partial}{\partial \theta} [\alpha_H \log \theta + \alpha_T \log(1-\theta)]$$

$$= \frac{\alpha_H}{\theta} - \frac{\alpha_T}{1-\theta} \quad \left[\text{assuming } \theta \neq 0 \right. \\ \left. \theta \neq 1 \right]$$

$$= \frac{\alpha_H - \alpha_H \theta - \alpha_T \theta}{\theta(1-\theta)} = 0$$

$$\Rightarrow \hat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}$$

② Sufficient Statistic

$\theta(X)$ is a sufficient statistic iff

$$\sum_{i \in D_1} \phi(x_i) = \sum_{i \in D_2} \phi(x_i) \Rightarrow L(\theta; D_1) = L(\theta; D_2)$$

Same statistics

Datasets appear "equivalent" to likelihood fn.

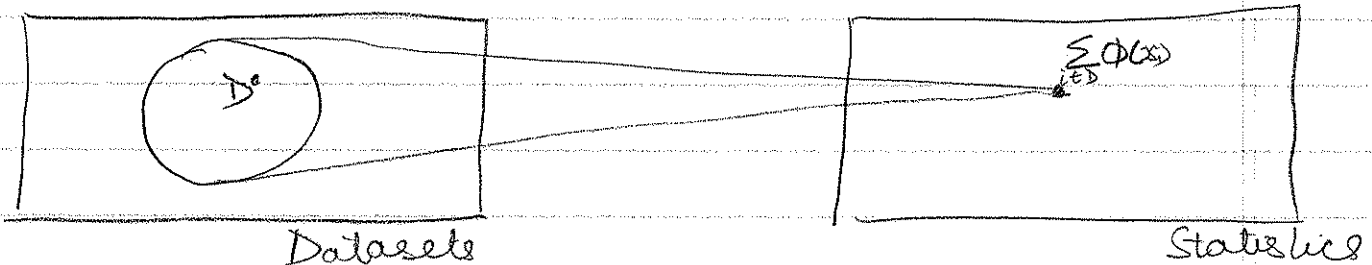
e.g. $D_1 = \{1, 1, 1, 0, 0, 0\}$

$D_2 = \{1, 0, 1, 0, 1, 0\}$

$\alpha_H = \alpha_T = 3$

$\alpha_H = \alpha_T = 3$

$\Rightarrow \phi(Y) = [Y=1]$



3

④ Why MLE? MLE is OPT if model class is correct!

→ Note: Strong Statement. Strong Assumption.

Consider, Y discrete (but argument generalizes)

$$\frac{1}{N} LL(\theta) = \frac{1}{N} \sum_{i=1}^N \log P(Y=y_i | \theta)$$

$$= \frac{1}{N} \left[\underbrace{\#(Y=1)}_{\text{count in dataset } D} \cdot \log P(Y=1 | \theta) + \#(Y=2) \cdot \log P(Y=2 | \theta) + \dots \right]$$

[∴ Counting argument]

As data becomes infinite,

$$\lim_{N \rightarrow \infty} \frac{\#(Y=1)}{N} = P(Y=1 | \theta^*)$$

↖ True parameter or unknown "reality"

Let's use shorthand $\left\{ \begin{array}{l} P^*(y) = P(Y=y | \theta^*) \\ P_\theta(y) = P(Y=y | \theta) \end{array} \right\}$

Now $\frac{1}{N} LL(\theta)$ as $N \rightarrow \infty$

$$= \sum_{y=1}^k P^*(y) \log P_\theta(y)$$

$$= \sum_{y=1}^k P^*(y) \log \left[\frac{P_\theta(y)}{P^*(y)} \cdot \frac{P^*(y)}{P_\theta(y)} \right]$$

$$= \underbrace{\sum_y P^*(y) \log P^*(y)}_{\text{entropy}} - \underbrace{\sum_y P^*(y) \log \frac{P^*(y)}{P_\theta(y)}}_{\text{KL-Divergence}}$$

$$\frac{1}{N} LL(\theta) = \underbrace{-H(p^*)}_{\text{entropy / uncertainty in } p^*} - \underbrace{KL(p^* \parallel P_0)}_{\text{How far is } P_0 \text{ from } p^*} \quad (\text{as } N \rightarrow \infty)$$

entropy /
uncertainty
in p^*
 \equiv constant
wrt θ

How far is P_0 from
 p^*

$$\rightarrow \underset{\theta}{\operatorname{argmax}} LL(\theta) = \underset{\theta}{\operatorname{argmin}} KL(p^* \parallel P_0)$$

Very cool!

POWERFUL RESULT

- We did not specify $P(Y=y|\theta)$ or the "Model Class"
- Result valid for any model!

Concave

- Inf dat
- We must know the "true" model P_0 , which we usually don't (e.g. Is life Gaussian?)

⑤ MAP + Bayesian Estimation

Key intuition: Let's think of θ as a random quantity & apply Bayes Rule

$$P(\theta | D) = \frac{P(D | \theta) P(\theta)}{P(D)}$$

Note: In Frequentist statistics, θ is unknown but not a random quantity

[e.g. there is only 1 world with 1 Nadal $\Rightarrow \theta$ is fixed] so can't talk about $P(\theta)$, but in Bayesian Stats:

$P(\theta)$ = Prior Belief
= what do we believe about θ without any data

In our running example $\theta \in [0, 1]$

\Rightarrow need a continuous distribution / density function

\Rightarrow a distribution over parameter of another distribution

\rightarrow Meet the Beta distribution

$$P(\theta | B_H, B_T) = \frac{\theta^{B_H-1} (1-\theta)^{B_T-1}}{\text{constant}}$$

hyper-parameters: parameters of distribution over parameter (θ)

Important Facts:

$$\rightarrow \text{constant} = \int_0^1 \theta^{\beta_H-1} (1-\theta)^{\beta_T-1} d\theta$$

Because pdf integrates to 1

$\theta \rightarrow$ mode of pdf

$$\hat{\theta}_{\text{mode}} = \frac{\beta_H - 1}{\beta_H + \beta_T - 2}$$



Maximum A Posteriori (MAP) Estimation

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\text{argmax}} P(\theta | D) \quad \left[\text{Let's pick the one value we "believe" in the most} \right]$$

$$= \underset{\theta}{\text{argmax}} \frac{P(D|\theta) P(\theta)}{P(D)}$$

constant w.r.t θ

$$= \underset{\theta}{\text{argmax}} P(D|\theta) P(\theta)$$

Special Case: $P(\theta) = \text{constant}$ [all θ s are equally likely]

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\text{argmax}} P(D|\theta)$$

$= \hat{\theta}_{\text{MLE}}$ [very nice, so frequentists are just Bayesians with no priors!]

(5)

In our setup,

$$P(\theta|D) \propto \theta^{\alpha_H} (1-\theta)^{\alpha_T} \times \frac{\theta^{B_H-1} (1-\theta)^{B_T-1}}{\text{constant}}$$

$$\propto \theta^{\alpha_H+B_H-1} (1-\theta)^{\alpha_T+B_T-1}$$

$$\propto \text{Beta}(\alpha_H+B_H, \alpha_T+B_T)$$

Very Nice! Conjugate Priors make math easy!

⇒ $\hat{\theta}_{\text{MAP}} = \text{mode of posterior Beta}$

$$= \frac{\alpha_H + B_H - 1}{\alpha_H + B_H + \alpha_T + B_T - 2}$$

So basically B_H, B_T act as "pseudo-flips"
→ experiments / data not contained in our dataset

Special Case: $\alpha_H = B_H = 1$

$$\hat{\theta}_{\text{MAP}} = \frac{\alpha_H + 1 - 1}{\alpha_H + 1 + \alpha_T + 1 - 2} = \hat{\theta}_{\text{MLE}}$$

Why? ∵ $P(\theta) \propto \theta^{1-1} (1-\theta)^{1-1} = \text{constant} / \text{uniform prior}$

