بسم الله الرحمن الرحيم

الكورس هيبقا عنيف .. ربنا ييسر الأمور ... هومورك رقم 6 نزل يا بيبيه ... خلص 5 انجز ... الدكتور عارف انك مش بتقرا ... اقرا يا بيبيه .... طيب الدكتور بيقول الحاجه الي عاوزين نركز عليها هي ال Max likelihood estimation و ال Max A-posterior estimation .. عشان نعمل كدا .. احنا هنبدأ نفكر في ال likelihood function .. ال setting بتاعها ان عندك N data samples .. دول كلهم sampled from random variable X ... تعال نبص علي أول حاله .. الإكس ديسكريت ... الضرب هنا حصل عشان هنا انت حاطط شرط iid ... تاني حاله .. ال likelihood function بقت فانكشن في الدينستي .. في الحالتين .. ال likelihood بتقولنا how likely it is one sample to another .. فلما بنتكلم علي continuous random variable انت هنا مش بتضرب احتمالات فبعض .. عادي .. انت بس مهتم بال relative likelihood مابين ال different settings of the parameter ...

## Likelihood Function — One R.V.

- Suppose we have N **samples** $D = \{x^{(1)}, x^{(2)}, \ldots, x^{(N)}\}$ from a **random variable** X

- The **likelihood** function: *prob. mass. function*   In both cases (discrete / continuous), the **likelihood** tells us how likely one sample is relative to another (iid)
  - <u>Case 1</u>: X is **discrete** with *pmf* $p(x|\theta)$
    $$L(\theta) = p(x^{(1)}|\theta)\, p(x^{(2)}|\theta) \ldots p(x^{(N)}|\theta)$$
  - <u>Case 2</u>: X is **continuous** with *pdf* $f(x|\theta)$
    $$L(\theta) = f(x^{(1)}|\theta)\, f(x^{(2)}|\theta) \ldots f(x^{(N)}|\theta)$$

- The **log**-**likelihood** function:
  - <u>Case 1</u>: X is **discrete** with *pmf* $p(x|\theta)$
    $$\ell(\theta) = \log p(x^{(1)}|\theta) + \ldots + \log p(x^{(N)}|\theta)$$
  - <u>Case 2</u>: X is **continuous** with *pdf* $f(x|\theta)$
    $$\ell(\theta) = \log f(x^{(1)}|\theta) + \ldots + \log f(x^{(N)}|\theta)$$

طيب احنا هنتكلم علي ال log likelihood function .. خد اللوج بس .. الكلام ده كلو ل One random variable .. اللي احنا كنا بنتكلم فيه طول الترم اصلا هو لما بيكون عندنا 2 راندم فاريبلز .. random variables x and y .. فاحنا عندنا x and y pairs .. سحبناهم من x and y .. اللي هنتكلم عنو انهارده هو ال joint likelihood function .. ده اللي هيطلعلنا لما نتكلم علي naïve Bayes ... عندنا حالتين برضو .. يا دسكريت يا كونتنوس ..

## Likelihood Function — Two R.V.s

- Suppose we have N **samples** $D = \{(x^{(1)}, y^{(1)}), \ldots, (x^{(N)}, y^{(N)})\}$ from a pair of **random variables** X, Y

- The **joint likelihood** function:
  - <u>Case 1</u>: X and Y are **discrete** with *pmf* $p(x,y|\theta)$
    $$L(\theta) = p(x^{(1)}, y^{(1)}|\theta) \ldots p(x^{(N)}, y^{(N)}|\theta)$$
  - <u>Case 2</u>: X and Y are **continuous** with *pdf* $f(x,y|\theta)$
    $$L(\theta) = f(x^{(1)}, y^{(1)}|\theta) \ldots f(x^{(N)}, y^{(N)}|\theta)$$
  
  Mixed discrete/continuous!
  
  - <u>Case 3</u>: Y is **discrete** with *pmf* $p(y|\beta)$ and X is **continuous** with *pdf* $f(x|y,\alpha)$
    $$L(\alpha, \beta) = f(x^{(1)}| y^{(1)}, \alpha)\, p(y^{(1)}|\beta) \ldots f(x^{(N)}| y^{(N)}, \alpha)\, p(y^{(N)}|\beta)$$
  - <u>Case 4</u>: Y is **continuous** with *pdf* $f(y|\beta)$ and X is **discrete** with *pmf* $p(x|y,\alpha)$
    $$L(\alpha, \beta) = p(x^{(1)}| y^{(1)}, \alpha)\, f(y^{(1)}|\beta) \ldots p(x^{(N)}| y^{(N)}, \alpha)\, f(y^{(N)}|\beta)$$

في حالتين كمان اللي هم Mixed discrete/continuous ... طيب افترض ان عندنا شوية داتا .. عاوزين نفكر في ال MLE شويه .. كريفيو يعني وكدا ..
احنا قلنا ان ال maximizing likelihood هدفو انو بيقول انا عندي finite amount of probability mass و كل اللي هو بيحاول يعمل انو ي
allocate much probability للحاجات اللي احنا لاحظناها عن الحاجات اللي احنا ملاحظنهاش ... فاحنا عندنا recipe للكلوزد فورم MLE ..

# Recipe for Closed-form MLE

1. Assume data was generated i.i.d. from some model
   (i.e. write the generative story)
   $x^{(i)} \sim p(x|\theta)$ ←

2. Write log-likelihood
   $\ell(\theta) = \log p(x^{(1)}|\theta) + \dots + \log p(x^{(N)}|\theta)$ ←

3. Compute partial derivatives (i.e. gradient)
   $\partial\ell(\theta)/\partial\theta_1 = \dots$
   $\partial\ell(\theta)/\partial\theta_2 = \dots$
   $\dots$
   $\partial\ell(\theta)/\partial\theta_M = \dots$

4. Set derivatives to zero and solve for $\theta$
   $\partial\ell(\theta)/\partial\theta_m = 0$ for all $m \in \{1, \dots, M\}$
   $\theta^{MLE} =$ solution to system of M equations and M variables

5. Compute the second derivative and check that $\ell(\theta)$ is concave down
   at $\theta^{MLE}$

poll: nlcourse.org

# MLE

| Question: | Answer: |
|---|---|
| Assume we have N samples $x^{(1)}$, $x^{(2)}, \dots, x^{(N)}$ drawn from a Bernoulli($\phi$). | A. $l(\phi) = N_1 \log(\phi) + N_0 (1 - \log(\phi))$ |
| | B. $l(\phi) = N_1 \log(\phi) + N_0 \log(1\text{-}\phi)$ |
| | C. $l(\phi) = \log(\phi)^{N_1} + (1 - \log(\phi))^{N_0}$ |
| What is the **log-likelihood** of the data $\ell(\phi)$? | D. $l(\phi) = \log(\phi)^{N_1} + \log(1\text{-}\phi)^{N_0}$ |
| | E. $l(\phi) = N_0 \log(\phi) + N_1 (1 - \log(\phi))$ |
| | F. $l(\phi) = N_0 \log(\phi) + N_1 \log(1\text{-}\phi)$ |
| Assume $N_1 = $ # of $(x^{(i)} = 1)$ | G. $l(\phi) = \log(\phi)^{N_0} + (1 - \log(\phi))^{N_1}$ |
| $N_0 = $ # of $(x^{(i)} = 0)$ | H. $l(\phi) = \log(\phi)^{N_0} + \log(1\text{-}\phi)^{N_1}$ |
| | I. $l(\phi) = $ the most likely answer |

Calamity

MLE of Bernoulli

① Model $x^{(i)} \sim Bernoulli(\phi)$

② Log-likelihood: $D = \{x^{(1)}, \dots, x^{(N)}\}$
$= \begin{cases} \phi & \text{if } x^{(i)} = 1 \\ 1-\phi & \text{if } x^{(i)} = 0 \end{cases} = \phi^{x^{(i)}}(1-\phi)^{(1-x^{(i)})}$

$\ell(\phi) = \log p(D|\phi)$

$= \log \prod_{i=1}^{N} p(x^{(i)}|\phi)$

$= \log \prod_{i=1}^{N} \phi^{x^{(i)}}(1-\phi)^{(1-x^{(i)})}$

$N_1 = \#(x^{(i)} = 1)$
$N_0 = \#(x^{(i)} = 0)$

$= \log [\phi^{N_1}(1-\phi)^{N_0}]$

$= N_1 \log\phi + N_0 \log(1-\phi)$     $\log p(D|\phi)$

# MLE

**Question:**

Assume we have N samples $x^{(1)}$, $x^{(2)}$, ..., $x^{(N)}$ drawn from a Bernoulli($\phi$).

What is the **derivative** of the log-likelihood $\partial \ell(\theta)/\partial \theta$?

Assume $N_1$ = # of $(x^{(i)} = 1)$

$\quad\quad N_0$ = # of $(x^{(i)} = 0)$

**Answer:**

~~calamity~~

A. ~~$\partial \ell(\theta)/\partial \theta = \phi^{N_1}(1-\phi)^{N_0}$~~

B. $\partial \ell(\theta)/\partial \theta = \phi / N_1 + (1 - \phi) / N_0$

C. $\partial \ell(\theta)/\partial \theta = N_1 / \phi - N_0 / (1 - \phi)$

D. $\partial \ell(\theta)/\partial \theta = \log(\phi) / N_1 + \log(1 - \phi) / N_0$

E. $\partial \ell(\theta)/\partial \theta = N_1 / \log(\phi) + N_0 / \log(1 - \phi)$

② Derivative

$$\frac{\partial \ell(\phi)}{\partial \phi} = \frac{\partial}{\partial \phi}\left[ N_1 \log \phi + N_0 \log(1-\phi) \right]$$

$$= \frac{N_1}{\phi} - \frac{N_0}{1-\phi}$$

③ Set to zero and solve

$$\frac{N_1}{\phi} - \frac{N_0}{1-\phi} = 0 \implies \phi^{MLE} = \frac{N_1}{N_1 + N_0} = \frac{N_1}{N}$$

# MLE vs. MAP

Suppose we have data $\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$

**Principle of Maximum Likelihood Estimation:**
Choose the parameters that maximize the likelihood of the data.

$$\boldsymbol{\theta}^{\text{MLE}} = \underset{\boldsymbol{\theta}}{\text{argmax}} \prod_{i=1}^{N} p(\mathbf{x}^{(i)}|\boldsymbol{\theta})$$

Maximum Likelihood Estimate (MLE)

**Principle of Maximum a posteriori (MAP) Estimation:**
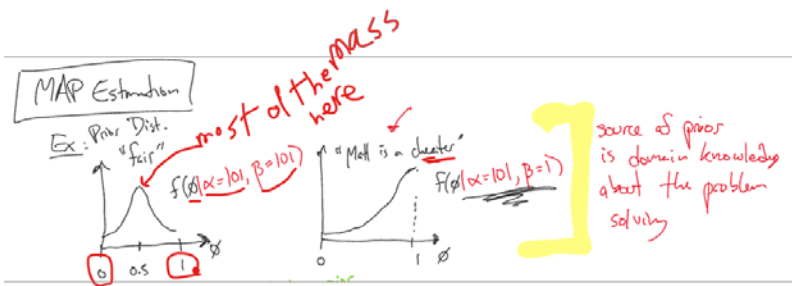Choose the parameters that maximize the posterior of the parameters given the data.

likelihood $\quad$ Prior

$$\boldsymbol{\theta}^{\text{MAP}} = \underset{\boldsymbol{\theta}}{\text{argmax}} \prod_{i=1}^{N} p(\mathbf{x}^{(i)}|\boldsymbol{\theta})p(\boldsymbol{\theta})$$

pdf

Maximum a posteriori (MAP) estimate

الدكتور بيقول ان ال MLE بيحاول ي (maximize p(D|phi) انما ال Max a posteriori بيحاول ي (Maximize p(phi|D) ...

## MAP Estimation

Ex: Prior Dist.

"fair"

most of the mass here

$f(\phi | \alpha=101, \beta=101)$

$0 \quad 0.5 \quad 1 \quad \phi$

"Matt is a cheater"

$f(\phi | \alpha=101, \beta=1)$

$0 \quad 1 \quad \phi$

source of prior is domain knowledge about the problem solving

---

$0 \quad 0.5 \quad 1 \qquad 0 \qquad 1 \quad x$

### General MAP: → ①

MLE: $p(D | \theta)$    likelihood    prior

MAP: $\underbrace{p(\theta | D)}_{\text{posterior}} = \dfrac{p(D|\theta) \, p(\theta)}{p(D)}$  → not a function of $\theta$

$\underbrace{\qquad}_{\text{Bayes Rule}}$

$\theta^{MAP} = \underset{\theta}{\arg\max} \; p(\theta | D)$ ← Posterior

$= \underset{\theta}{\arg\max} \; \log p(\theta | D)$

$= \underset{\theta}{\arg\max} \; \log \left( \dfrac{p(D|\theta) p(\theta)}{\boxed{p(D)}} \right)$ → doesn't affect argmax ("remove it")

$= \underset{\theta}{\arg\max} \; \underbrace{\log \big( p(D|\theta) p(\theta) \big)}_{\ell_{MAP}(\theta)}$

$\checkmark^{MAP}$

---

## MAP of Beta-Bernoulli

Beta Dist:

$f(\phi | \alpha, \beta) = \dfrac{1}{B(\alpha, \beta)} \phi^{(\alpha - 1)} (1 - \phi)^{(\beta - 1)}$

① Model:   $\phi \sim \text{Beta}(\alpha, \beta)$

$x^{(1)} \sim \text{Bernoulli}(\phi)$

$x^{(2)} \sim \text{Bernoulli}(\phi)$

$\vdots$

$x^{(N)} \sim \text{Bernoulli}(\phi)$

parameters    hyperparameters

② Log-likelihood

$\ell_{MAP}(\phi) = \log \Big[ p(D | \phi) \, f(\phi | \alpha, \beta) \Big]$

$= \log \Big[ \phi^{N_1} (1 - \phi)^{N_0} \left( \dfrac{1}{B(\alpha, \beta)} \phi^{(\alpha-1)} (1 - \phi)^{(\beta-1)} \right) \Big]$

$= \log \Big[ \phi^{(N_1 + \alpha - 1)} (1 - \phi)^{(N_0 + \beta - 1)} \dfrac{1}{B(\alpha, \beta)} \Big]$

$= (N_1 + \alpha - 1) \log(\phi) + (N_0 + \beta - 1) \log(1 - \phi) - \log(B(\alpha, \beta))$

$= N_1' \log(\phi) + N_0' \log(1 - \phi) - \log(B(\alpha, \beta))$

③ Derivative

$$\frac{d\ell_{MAP}(\emptyset)}{\partial \emptyset} = \frac{N_1'}{\emptyset} - \frac{N_0'}{1-\emptyset}$$

④ Set to zero and solve

$$\emptyset^{MAP} = \frac{N_1'}{N_1' + N_0'} = \frac{N_1 + \alpha - 1}{N_1 + \alpha - 1 + N_0 + \beta - 1}$$

Real life !~

Prior dist. fair

لو انت دخلت prior dist. fair لبقاء ال 8 و ال 2 طلعو كانوا noise

لو ر فلّت prior MLE

متوقع ان ال داتا ما تش كتّر

same for (101, 1)

عندك 10 coin flips

قبل ما نلعب و بعد ما دد فلفشتوفور نلعب ولا لأ؟

$N_1 + N_0$

Ex#1 : Suppose $D = \{8H, 2T\}$

$\emptyset^{MLE} = \frac{8}{10} = 0.8$

Prior dist. → Now if $\emptyset \sim Beta(\alpha = 101, \beta = 101)$

pseudo counts → to heads & tails MLE

fair

cheat

$\emptyset^{MAP} = \frac{8 + 101 - 1}{8 + 101 - 1 + 2 + 101 - 1} = \frac{108}{108 + 102} \approx 0.5$

Now if $\emptyset \sim Beta(\alpha = 101, \beta = 1)$

$\emptyset^{MAP_2} = \frac{108}{108 + 2} \approx 1.0$

Ex#2 : $D = \{108H, 102T\}$

$\emptyset^{MLE} = \frac{108}{108 + 102} = \emptyset^{MAP_1}$

Prior Beta

و هنا يبقى ده هو ال Prior counts pseudo

يعني إن كانت لا تعتمد القيم على ال experiment

$\alpha = 100$ ('heads') $\beta = 100$ ('tails')

الدكتور اتكلم علي حتة ال pseudo counts وقال ان لما انت بتلاحظ عد كبير من الداتا .. ال MLE هتبقا نفس القيمه بتاعت ال MAP رقم 1 من الصوره اللي فوق ... طيب احنا كل اللي بنحاول نعملو هو function approximation و ال MLE بيدينا view مختلف شويه للتعليم .. للي بنحاول نعملو هو اننا ن Max likelihood of the data بس في الآخر الدنيا بتوصل برضو لل function approximation ...

هنتكلم ناو علي Naïve Bayes .. هو مجرد decision function ...

# Naïve Bayes

- Why are we talking about Naïve Bayes?
  - It's **just another decision function** that fits into our "big picture" recipe from last time
  - But it's our first **example of a Bayesian Network** and provides a *clearer* picture of **probabilistic learning**
  - Just like the other Bayes Nets we'll see, it **admits a closed form solution** for MLE and MAP
  - So learning is **extremely efficient** (just counting)

---

MLE for Naive Bayes

Data: $\begin{cases} y \in \{H,T\} \\ y \in \{0,1\} \\ \vec{x} \in \{0,1\}^M \end{cases}$   where real = 1 and false = 0
where $M = \#$ words in vocab.

① Model:
$$y \sim \text{Bernoulli}(\phi) = p(y|\phi)$$
$$X_1 \sim \text{Bernoulli}(\theta_{y,1}) = p(x_1|y,\theta)$$
$$X_2 \sim \text{Bernoulli}(\theta_{y,2}) = p(x_2|y,\theta)$$
$$\vdots$$
$$X_M \sim \text{Bernoulli}(\theta_{y,M}) = p(x_M|y,\theta)$$

$\phi \in [0,1]$    $\in [0,1]$

$\theta = \begin{bmatrix} \theta_{H1} & \theta_{H2} & \cdots & \theta_{HM} \\ \theta_{T1} & \theta_{T2} & \cdots & \theta_{TM} \end{bmatrix}$ ← red Corns ← blue Corn

Tail

$$P(x_1, x_2, \ldots, x_M, y | \phi, \theta) = p(y|\phi) p(x_1|y,\theta) p(x_2|y,\theta) \cdots p(x_M|y,\theta)$$

words

Real or fake

$$= p(y|\phi) \prod_{M=1}^{M} p(x_M|y, \theta_{H,M}, \theta_{T,M})$$

From Bernoulli: $= \phi^y (1-\phi)^{(1-y)} \prod_{M=1}^{M} \theta_{y,M}^{X_M} (1-\theta_{y,M})^{(1-X_M)}$

Def: two r.v.s $X, Y$ are Conditionally independent given r.v. $Z$ written $X \perp\!\!\!\perp Y | Z$ iff
$$p(X,Y|Z) = p(X|Z)p(Y|Z)$$

---

Naive Bayes Assumption

$$p(\vec{x}|y) = \prod_{M=1}^{M} p(x_M|y)$$

$\implies X_q$ and $X_r$ are cond. indep. given $y$