# Probabilistic modeling, MLE and MAP Estimates

Recall the ML setup: $(x, y) \sim P$

If we knew $P(X, Y)$ or even just $P(Y|X)$, we could compute Bayes optimal classifier

For classification

more generally

$$h(x) = \underset{y \in C}{\text{argmax}} \; P(Y = y \mid X = x)$$

$$h(x) = \underset{\hat{y}}{\text{argmin}} \; E_{y|x}[\ell(\hat{y}, y)]$$

Generative: $P(X|Y) \, P(Y)$
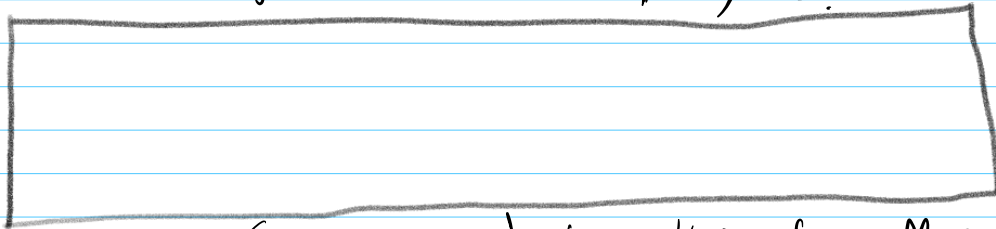
## Probabilistic modeling: Estimate $P(X, y)$

Discriminative (or $P(Y|X)$ directly) and use it instead

Estimating Bernoilli R.V.: Yearly rain/no Rain Data

$$D = \{R, N, N, N, R, N, N\}$$

$R = $ Rain    $N = $ No Rain

$X = \{\}$
$C = \{R, N\}$

What would your estimate to $P(Y)$ be?

Can we derive this formally?

Assume events are "Independent and Identically distributed"
i.i.d.

Parameter: $p = P(Y = R)$    $n_R = $ # Rainy days
$n_N = $ # no rain days

$$\hat{p} = \underset{p \in [0,1]}{\text{argmax}} \binom{n_R + n_N}{n_R} p^{n_R} (1-p)^{n_N} \quad \textcircled{1}$$

$$= \underset{p \in [0,1]}{\text{argmax}} \log\left(\binom{n_R + n_N}{n_R}\right) + n_R \log p + n_N \log(1-p) \quad \textcircled{2}$$

$$= \underset{p \in [0,1]}{\text{argmax}} \quad n_R \log p + n_N \log(1-p) \quad \textcircled{3}$$

$$\textcircled{4} \implies \frac{n_R}{p} - \frac{n_N}{1-\hat{p}} = 0 \implies \hat{p} = \frac{n_R}{n_R + n_N}$$

1. Parameterize $P(\text{Data})$ by some family of parameters $P_\theta$ s.t. $\theta \in \Theta$

2. Estimate $P(x,y)$ (or $P(y|x)$) by estimating $\theta \in \Theta$ from Data

M LE : Maximum Likelihood Estimate

$$\theta_{MLE} = \underset{\theta \in \Theta}{\text{argmax}} \quad P_\theta(\text{Data})$$

pick $\theta$ that maximizes likelihood of Data observation

1 Often referred to as frequentist view
2 when $\theta^* \in \Theta$ generates data, $\theta_{MLE} \to \theta^*$ (typically)

Eg 2: Data : $D = \{176, 177, 169, 168, \dots\}$
Heights of Adult Male/Female

heights are normally distributed with mean $\mu$ and variance $\sigma^2$

$$\Theta = \{(\mu, \sigma^2) : \mu \in \mathbb{R}, \sigma^2 \in \mathbb{R}^+\}$$

① $\left(\widehat{\mu}_{MLE}, \widehat{\sigma}^2_{MLE}\right) = \underset{\mu, \sigma^2}{argmax} \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2}\frac{(x_i-\mu)^2}{\sigma^2}\right)$
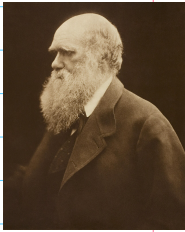
② $\left(\widehat{\mu}_{MLE}, \widehat{\sigma}^2_{MLE}\right) = \underset{\mu, \sigma^2}{argmax} -\sum_{i=1}^{n}\left(\frac{(x_i-\mu)^2}{2\sigma^2} + \log\left(\frac{1}{\sigma}\right)\right)$

③ $\frac{\sum_{i=1}^{n}(x_i - \widehat{\mu}_{MLE})}{\sigma^2} = 0$ , $\sum_{i=1}^{n}\left(\frac{(x_i - \widehat{\mu}_{MLE})^2}{\widehat{\sigma}^3_{MLE}} - \frac{1}{\widehat{\sigma}_{MLE}}\right) = 0$
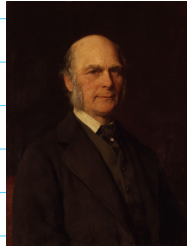
④ $\left(\widehat{\mu}_{MLE}, \widehat{\sigma}^2_{MLE}\right) = \left(\frac{1}{n}\sum_{i=1}^{n} x_i , \frac{1}{n}\sum_{i=1}^{n}(x_i - \widehat{\mu}_{MLE})^2\right)$

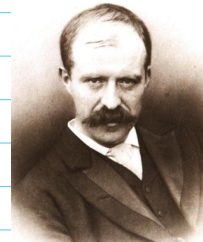Sample mean        Sample variance

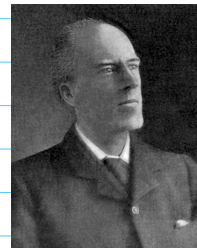Gaussian Mixture Model :


Charles Darwin


Francis Galton


Raphael Weldon


Karl Pearson

Vs

Evolution via Natural selection

Evolution discontinuous Sudden

small variations, gradual evolution

$\Pi_1 = \frac{1}{2}$

$\Pi_2 = \frac{1}{2}$




Plot #012   Data: Pearson's crab   Components: Normal

$\Theta = \begin{cases} \overline{\Pi} = (\Pi_1, \dots, \Pi_K) & \text{distribution over } K \text{ items} \\ \mu_1, \dots, \mu_K \in \mathbb{R}^d & K - means \\ \Sigma_1, \dots, \Sigma_K \in \mathbb{R}^{d\times d} & K - covariances \end{cases}$

Normal pdf
↑

$\widehat{\Theta}_{MLE} = \underset{\Pi, \mu_1 \dots \mu_K, \Sigma_1 \dots \Sigma_K}{argmax} \quad P_\Theta(D) := \prod_{i=1}^{n}\left(\sum_{k=1}^{K} \Pi_k \, p(x_i \mid \mu_k, \Sigma_k)\right)$

EM Algo for GMM aims to solve above

# MLE does not capture prior knowledge

Eg Rain, No Rain.

Say we had prior info. that at similar locations typically we have seen Rain on 30 out of 100 days, how do we use this?

Heuristic: $p = P(Y = Rain) = \dfrac{n_R + 30}{n_R + n_N + 100}$
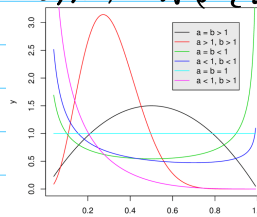
## Maximum Aposteriori Estimator: MAP

Model is an abstraction that captures our belief, we update our belief based on Data.

$\theta$ is a Random variable

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

$$= \underset{\theta \in \Theta}{argmax} \; \underbrace{\log(P(D|\theta))}_{\text{log likelihood}} + \underbrace{\log P(\theta)}_{\text{log prior}}$$

Eg: For Bernoulli distribution we can use Beta prior

$$P(\theta) = \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{B(\alpha, \beta)}$$



① $\hat{\theta}_{MAP} = \underset{\theta}{argmax} \; \log P(D|\theta) + \log P(\theta)$

② $= \underset{\theta}{argmax} \; n_R \log \theta + n_N \log(1-\theta) +$

$\qquad (\alpha-1) \log \theta + (\beta-1) \log(1-\theta) - \log B(\alpha,\beta)$

③ $\hat{\theta}_{MAP} = \dfrac{n_R + \alpha - 1}{n_R + n_R + \alpha + \beta - 2}$

$\alpha - 1$ Rains
$\beta - 1$ No Rains

Often MAP is referred to as Bayesian view

There is Bayesian and there is BAYESIAN

True Bayesian: "There is no model, all you are estimating is $Y$"

$$P(Y|X, Data) = \int_\theta P(Y, \theta | X, Data) \, d\theta$$

$$= \int_\theta P(Y | \theta, X, Data) \, P(\theta | Data) \, d\theta$$