

Intro about Bayesian Machine Learning

• Setting: $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$

• Up to now, the ML methods we used can be seen as based on max likelihood

$$\hat{\theta}_{ML} = \arg \max_{\theta} p(\mathcal{D} | \theta)$$

or max a posteriori

$$\hat{\theta}_{MAP} = \arg \max_{\theta} p(\theta | \mathcal{D})$$

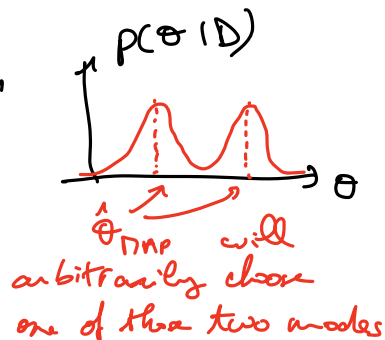
$$= \arg \max_{\theta} \underbrace{\log p(\mathcal{D} | \theta)}_{\text{log likelihood}} + \underbrace{\log p(\theta)}_{\text{regularization term}}$$

because

$$p(\theta | \mathcal{D}) = \frac{p(\mathcal{D} | \theta) p(\theta)}{p(\mathcal{D})} \leftarrow \text{constant}$$

• We were looking for a single model $\hat{\theta}$

• pb: - no level of confidence in our choice
- maybe there is no good "single model"



I] A simple coin estimation problem
with the ML, MAP, and bayesian views.
↳ max likelihood.

A] Max Likelihood approach

• We have a coin with $p(\text{heads}) = \theta$ unknown.

• We have drawn the coin N times, we got n_H heads and n_T tails.

• We want to estimate $p(\text{heads})$ based on n_H, n_T

• simplest approach is $\hat{\theta}_{ML} = \frac{n_H}{n_H + n_T} \leftarrow \text{max likelihood}$

- derivation of $\hat{\theta}_{ML}$:

- the likelihood is $P(\underbrace{n_H, n_T}_D | \theta) = \binom{n_H + n_T}{n_H} \theta^{n_H} (1-\theta)^{n_T}$
(binomial law)

$$\hat{\theta}_{ML} = \underset{\theta}{\operatorname{argmax}} \log P(n_H, n_T | \theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \underbrace{n_H \log \theta + n_T \log(1-\theta)}$$

derive and set derivative to zero

$$\frac{n_H}{\hat{\theta}} = \frac{n_T}{1-\hat{\theta}} \Rightarrow n_H - n_H \hat{\theta} = n_T \hat{\theta} \Rightarrow \hat{\theta}_{ML} = \frac{n_H}{n_H + n_T}$$

B) Max a Posteriori approach

Note: D is the pair (n_H, n_T)

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

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

→ likelihood

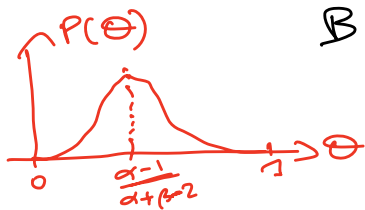
→ prior

→ evidence (cte)

We need a prior on θ .

Typically, we use a Beta distribution.

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



Note:

$$B(\alpha, \beta) = \frac{\Gamma(\alpha) \Gamma(\beta)}{\Gamma(\alpha + \beta)}$$

$$E[\theta] = \frac{\alpha}{\alpha + \beta}$$

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

$$= \underset{\theta}{\operatorname{argmax}} \text{cte} + n_H \log \theta + n_T \log(1-\theta) + (\alpha-1) \log \theta + (\beta-1) \log(1-\theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \underbrace{(n_H + \alpha - 1) \log \theta + (n_T + \beta - 1) \log(1-\theta)}$$

= log P(θ|D) + cte = log beta distribut. w.r. π, θ

$$\hat{\theta}_{MAP} = \frac{n_H + \alpha - 1}{n_H + n_T + \alpha + \beta - 2}$$

The Beta prior is said to be the "conjugate" prior of the binomial because

1) the posterior $P(\theta | D)$ is a binomial distribution and a beta distribution w.r.t. θ with parameters $\text{Beta}(\alpha + n_H, \beta + n_T)$

2) adding the prior has the same effect as observing additional draws.

Note that if $N \rightarrow \infty$ then $\hat{\theta}_{MAP} \rightarrow \hat{\theta}_{ML}$

c) Full Bayesian approach

- What is the probability that our next draw is a head?

$$p(\text{Head} | D) = \int_{\theta} P(\text{Head}, \theta | D) d\theta$$

(Law of total prob.)

$$= \int_{\theta} P(\text{Head} | \theta, D) \times P(\theta | D) d\theta$$

$$= \int_{\theta} \underbrace{P(\text{Head} | \theta)}_{\theta} \times P(\theta | D) d\theta$$

$$= \int_{\theta} \theta \cdot P(\theta | D) d\theta$$

$$= \mathbb{E}_{\theta \sim P(\theta | D)} [\theta | D]$$

↓
 $\text{Beta}(\alpha + n_H, \beta + n_T)$

$$P(\text{Heads} | D) = \frac{n_H + \alpha}{n_H + n_T + \alpha + \beta}$$

II) Bayesian view of Machine learning

- Assume $D = \{(x_i, y_i)\}_{i=1}^N$ (classification or regression)
- Let $X = (x_1, \dots, x_N)$ $Y = (y_1, \dots, y_N)$
- Assume X is fixed
- Assume the "nature" chose θ from $p(\theta)$
- Y is drawn from $P(Y | X, \theta) = P(Y | \theta)$

Let x_{N+1} be a new data point on which a prediction has to be made.

In max likelihood: $\hat{\theta}_{ML} = \arg\max_{\theta} P(Y | \theta)$

In the Bayesian setting, we don't want to choose a specific $\hat{\theta}$, we only want to predict y_{N+1}

$$P(y_{N+1} | X, Y, x_{N+1})$$

(predictive distribution)

$$= P(y_{N+1} | Y) = \int P(y_{N+1}, \theta | Y) d\theta$$

total law of prob

$$= \int P(y_{N+1} | \theta, Y) p(\theta | Y) d\theta$$

$$= \int P(y_{N+1} | \theta) p(\theta | Y) d\theta$$

$$\propto \int \underbrace{P(y_{N+1} | \theta)}_{\text{prediction}} \underbrace{P(Y | \theta)}_{\text{likelihood}} \underbrace{p(\theta)}_{\text{prior for } \theta} d\theta$$

Bayesian Model Averaging

III Approximation of the BAI integral.

x_{N+1} is constant as well as x , so no need to put them in the conditional.

$$\underbrace{p(y_{N+1} | \mathcal{Y})}_{\text{predictive proba distribution}} \propto \int p(y_{N+1} | \theta) p(\mathcal{Y} | \theta) p(\theta) d\theta$$

Naïve idea: draw randomly $\theta_1, \dots, \theta_n$ from $p(\theta)$
e.g. $p(\theta) = \mathcal{N}(0, \text{Id})$

$$p(y_{N+1} | \mathcal{Y}) \approx \frac{1}{Z} \sum_{k=1}^n p(y_{N+1} | \theta_k) \cdot \underbrace{p(\mathcal{Y} | \theta_k)}_{\text{eikd-ked}}$$

$$Z = \sum_{k=1}^n p(\mathcal{Y} | \theta_k)$$

Can't work because n would need to be huge for the sum to approximate the integral.

Better Idea: draw $\theta_1, \dots, \theta_n$ from $p(\theta | \mathcal{Y})$ the posterior.

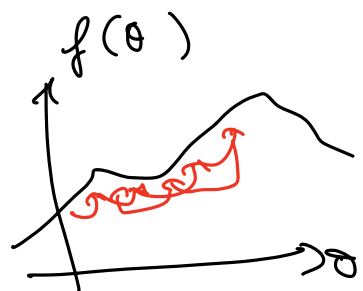
$$p(y_{N+1} | \mathcal{Y}) \approx \frac{1}{n} \sum_{k=1}^n p(y_{N+1} | \theta_k)$$

Much better estimator!

sampling from $p(\theta | \mathcal{Y})$ is hard.

A) Stochastic Langevin gradient for sampling from $p(\theta | \mathcal{Y})$

- let $f(\theta)$ be a density over θ .
- assume we want to sample $\theta \sim f$



- One way is to use Langevin stochastic gradient.

Repeat:

$$\theta \leftarrow \theta + \frac{\alpha}{2} \nabla \log f + \varepsilon_t$$

After "some" steps of the algorithm, θ is distributed according to $f(\theta)$

$$\begin{aligned} \sum_{t=0}^{\infty} \alpha_t &= \infty \\ \sum_{t=0}^{\infty} \alpha_t^2 &< \infty \\ \varepsilon_t &\sim \mathcal{N}(0, \alpha_t \cdot \text{Id}) \end{aligned}$$

- In our machine learning setting, we do a gradient descent on

$-\log p(\theta | \mathcal{Y})$ and we add noise on θ at each step.

- the successive θ s are kept as samples of $p(\theta | \mathcal{Y})$

B) Variational approximation

We will learn $q(\theta) = \mathcal{N}(\theta | \mu, \Sigma)$ which approximate $p(\theta | \mathcal{Y})$ as well as possible.

I want to minimize

$$\text{KL}(q(\theta) \parallel p(\theta | \mathcal{Y})) = \int q(\theta) \log \frac{q(\theta)}{p(\theta | \mathcal{Y})} d\theta$$

$$\frac{p(\theta | \mathcal{Y})}{q(\theta)} = \frac{p(\mathcal{Y} | \theta) p(\theta)}{q(\theta) \cdot p(\mathcal{Y})}$$

$$\begin{aligned}
KL(q(\theta) \parallel p(\theta | \mathcal{Y})) &= - \int q(\theta) \log \frac{p(\theta | \mathcal{Y})}{q(\theta)} d\theta \\
&= - \int q(\theta) \cdot [\log p(\mathcal{Y} | \theta) + \log \left(\frac{p(\theta)}{q(\theta)} \right) - \log p(\mathcal{Y})] \\
&= - \underbrace{\mathbb{E}_{\theta \sim q} [\log p(\mathcal{Y} | \theta)]}_{\text{negative log likelihood}} + \underbrace{KL(q(\theta), p(\theta))}_{\text{complexity term}} + \underbrace{\log p(\mathcal{Y})}_{\text{evidence}}
\end{aligned}$$

ELBO

To minimize $KL(q(\theta) \parallel p(\theta | \mathcal{Y}))$, I will use gradient descent with the reparameterization trick:

let $\Sigma = \begin{pmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_d^2 \end{pmatrix}$ recall $q(\theta) = \mathcal{N}(\theta | \mu, \Sigma)$

$$\mathbb{E}_{\theta \sim q} [f(\theta)] = \mathbb{E}_{\varepsilon \sim \mathcal{N}(0, \text{Id})} [f(\mu + \sqrt{\Sigma} \cdot \varepsilon)]$$

$$\begin{aligned}
KL(q(\theta) \parallel p(\theta | \mathcal{Y})) &= - \mathbb{E}_{\varepsilon \sim \mathcal{N}(0, \text{Id})} [\log p(\mathcal{Y} | \mu + \sqrt{\Sigma} \cdot \varepsilon)] \\
&\quad + KL(q(\theta), p(\theta)) + \text{cte.}
\end{aligned}$$

Learning algorithm:

$$\begin{aligned}
&\text{repeat} \\
&\quad \varepsilon \sim \mathcal{N}(0, \text{Id}) \\
&\quad \begin{pmatrix} \mu \\ \Sigma \end{pmatrix} \leftarrow \begin{pmatrix} \mu \\ \Sigma \end{pmatrix} - \nabla_{\begin{pmatrix} \mu \\ \Sigma \end{pmatrix}} \left\{ -\log p(\mathcal{Y} | \mu + \sqrt{\Sigma} \cdot \varepsilon) + KL(\mathcal{N}(\mu, \Sigma), p(\theta)) \right\}
\end{aligned}$$

To simplify, assume we set $\Sigma = \text{Id}$, and we only learn μ .
and $p(\theta) = \mathcal{N}(\theta, \text{Id})$. The algo becomes:

$$\left[\begin{array}{l} \text{repeat} \\ \quad \varepsilon \sim \mathcal{N}(\mathbf{0}, \text{Id}) \\ \quad \mu \leftarrow \mu - \nabla_{\mu} \left\{ -\log(P(\mathcal{Y} | \mu + \varepsilon)) + \lambda \|\mu\|^2 \right\} \end{array} \right]$$

to make a prediction

$$p(y_{N+1} | \mathcal{Y}) = \int p(y_{N+1} | \theta) \mathcal{N}(\theta | \mu, \Sigma) d\theta \\ \approx \frac{1}{M} \sum_{k=1}^M p(y_{N+1} | \theta_k)$$

$$\text{with } \theta_k \sim \mathcal{N}(\mu, \Sigma)$$

Note: Langevin SGD requires us to store M models $\theta_1 \dots \theta_M$,
but Variational methods only require
us to store μ, Σ which is lighter.

IV Bayesian linear regression

A). $D = \{(x_i, y_i)\}_{i=1}^N$, $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$

Assume there exists θ^* such that

$$y_i = \theta^{*T} x_i + \varepsilon_i \quad \text{where } \varepsilon_i \sim \mathcal{N}(0, 1)$$

$$\text{so } y_i \sim \mathcal{N}(\cdot | \theta^{*T} x_i, 1)$$

$$p(Y|\theta) = \prod_{i=1}^N p(y_i|\theta) = \prod_{i=1}^N \mathcal{N}(y_i | \theta^T x_i, 1)$$

. I define the prior $p(\theta) = \mathcal{N}(0, \sigma^2 \text{Id})$

$$p(\theta|Y) = \frac{p(Y|\theta) \cdot p(\theta)}{p(Y)} \quad \text{is gaussian.}$$

$$\begin{aligned} \log p(\theta|Y) &= \log p(Y|\theta) + \log p(\theta) + \text{cte} \\ &= -\frac{1}{2} \|Y - X\theta\|^2 - \frac{1}{2\sigma^2} \|\theta\|^2 + \text{cte} \\ &= -\frac{1}{2} \left(Y^T Y - 2\theta^T X^T Y + \theta^T X^T X \theta + \frac{1}{\sigma^2} \|\theta\|^2 \right) + \text{cte} \\ &= -\frac{1}{2} \theta^T \left(\underbrace{X^T X + \frac{1}{\sigma^2} \text{Id}}_{\Sigma^{-1}} \right) \theta + \underbrace{\theta^T X^T Y}_{\Sigma^{-1} \mu} + \text{cte} \end{aligned}$$

I know $p(\theta|Y)$ is gaussian.

$$\text{let } p(\theta|Y) = \mathcal{N}(\mu, \Sigma)$$

Now we have to identify μ and Σ .

$$\begin{aligned} \log p(\theta|Y) &= -\frac{1}{2} (\theta - \mu)^T \Sigma^{-1} (\theta - \mu) + \text{cte} \\ &= -\frac{1}{2} \theta^T \Sigma^{-1} \theta + \theta^T \Sigma^{-1} \mu + \text{cte} \end{aligned}$$

$$\Sigma = \left(\cancel{X^T X} + \frac{1}{\sigma^2} \text{Id} \right)^{-1}$$

$$\mu = \Sigma X^T Y = \left(\cancel{X^T X} + \frac{1}{\sigma^2} \text{Id} \right)^{-1} X^T Y$$

recall that this coincide with $\hat{\theta}_{\text{MAP}}$

$p(\theta|Y)$ is centered on $\hat{\theta}_{\text{MAP}}$

$$p(y_{N+1} | Y) = \int_{\theta} p(y_{N+1} | \theta) p(\theta | Y) d\theta$$

$$= \int_{\theta} \mathcal{N}(y_{N+1} | \theta^T x_{N+1}, 1) \mathcal{N}(\theta | \mu, \Sigma) d\theta$$

$$= \mathcal{N}(\cdot | \mu^T x_{N+1}, 1 + x_{N+1}^T \Sigma x_{N+1})$$

B) Non linear bayesian regression

like for the kernel trick,
 we use a mapping $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$
 we will learn a predictor $f_{\theta}(x) = \theta^T \phi(x)$
⏟
non linear in x

Assume $\theta \sim p(\theta) = \mathcal{N}(0, Id)$

then $\mathbb{E}_{\theta} [f_{\theta}(x)] = 0$

$$\begin{aligned} \text{Cov}_{\theta} [f_{\theta}(x), f_{\theta}(x')] &= \mathbb{E}_{\theta} [\phi(x)^T \theta \theta^T \phi(x')] \\ &= \phi(x)^T \underbrace{\mathbb{E}_{\theta} [\theta \theta^T]}_{Id} \phi(x') \\ &= \phi(x)^T \phi(x') \\ &= k(x, x') \end{aligned}$$

If $\theta \sim \mathcal{N}(0, Id)$

then $f_{\theta}(x)$ is gaussian

and the vector

$$(f_{\theta}(x_1), \dots, f_{\theta}(x_n)) \sim \mathcal{N}(0, K)$$

where $K_{ij} = k(x_i, x_j)$

This is a gaussian process.

g s y