

# Bayesian Deep Learning

Jan 9, 2021

NASH & IASD

Acknowledgement:

I've borrowed many ideas  
from Andrew Gordon Wilson, NYU.

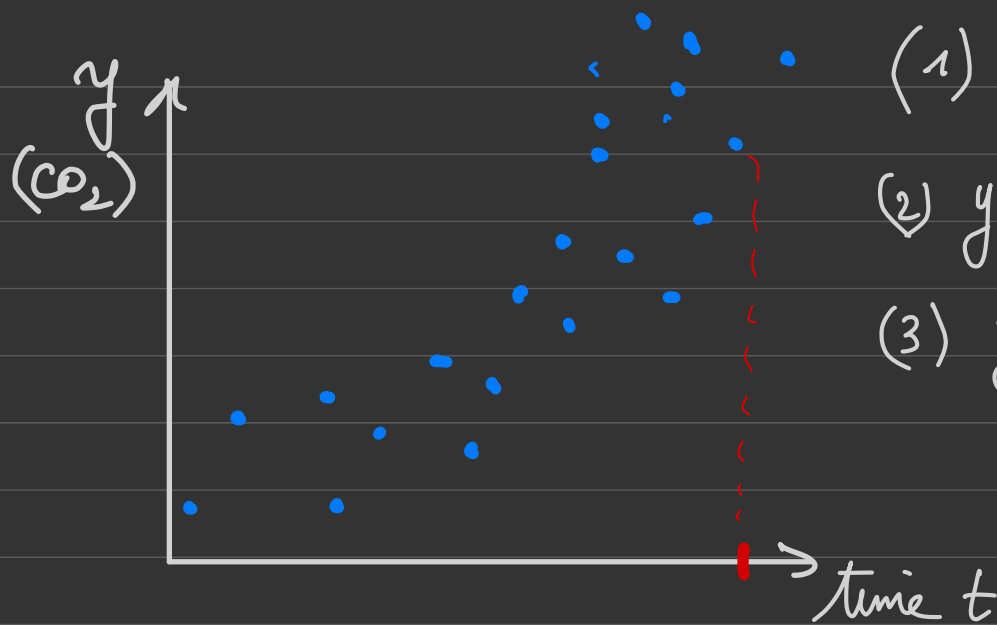
What comes to your mind when you hear

Bayesian inference!

probability distribution from data

update probs. given initial belief

Deep Learning



$$(1) y(t) = a_0 + a_1 t$$

$$(2) y(t) = \sum_{i=0}^3 a_i t^i$$

$$(3) y(t) = \sum_{i=0}^{10^4} a_i t^i$$

Uncertainty:

- o when gathering data : aleatoric uncertainty
  - irreducible uncertainty
- o choice of model / hypothesis
  - epistemic uncertainty
  - reducible uncertainty

$$f(x; w) = \sum_{j=0}^J w_j x^j$$

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1 \dots n}$$

$$\text{learn } \underline{w} = (w_0, w_1, \dots, w_J)$$

$$\mathcal{L}_{SE}(\underline{w}) = \sum_{i=1}^n \left( f(x_i; \underline{w}) - y_i \right)^2$$

$$\text{Minimize wrt } \underline{w}$$

$$\mathcal{L}_{AE}(\underline{w}) = \sum_{i=1}^n \left| f(x_i; \underline{w}) - y_i \right|$$

$$y(x) = f(x; \underline{w}) + \varepsilon(x) \quad \text{noise}$$

$$y(x) \sim N(f(x, \underline{w}), \sigma^2) \Rightarrow \varepsilon(x) = \varepsilon \sim N(0, \sigma^2)$$

likelihood  $\mathcal{D} = \{(x_i, y_i)\}_{i=1 \dots n}$

$$p(y_1, \dots, y_n | x_1, \dots, x_n, \underline{w})$$

$$= \prod_{i=1}^n p(y_i | x_i, \underline{w}) \quad (\text{by conditional iid})$$

$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2} (f(x_i, \underline{w}) - y_i)^2\right)$$

$$\log\text{-like} = \log\left[\left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n\right] - \underbrace{\frac{1}{2\sigma^2} \sum_{i=1}^n (f(x_i, \underline{w}) - y_i)^2}_{\text{Same as previous objective.}}$$

Same as previous objective.

▣ Gaussian noise  $\leadsto$  squared error.

▣ What about other distr. for noise

e.g. Laplace:  $p(w) = \frac{1}{2} e^{-|x|}$

heavier than Gaussian

# Bayesian modeling

Like.  $p(y | x, \underline{w}) = N(y; f(x; \underline{w}), \sigma^2)$

Prior.  $p(\underline{w}) = N(0, \sigma^2 I)$

$$\Leftrightarrow \forall j, p(w_j) = N(0, \sigma^2)$$

Sum rule  $p(a) = \sum_b p(a, b)$

Product rule  $p(a, b) = p(a) p(b|a)$   
 $= p(b) p(a|b)$

$\Rightarrow$  Posterior:  $p(\underline{w} | \mathcal{D}) = \frac{p(\underline{w}) p(\mathcal{D} | \underline{w})}{p(\mathcal{D})}$

*prior like.*  
} marginal likeli.  
evidence

Take log:

$$\log p(\underline{w} | \mathcal{D}) = \log p(\underline{w}) + \log p(\mathcal{D} | \underline{w}) - \log p(\mathcal{D})$$

*not def. on  $\underline{w}$*

$$\text{Penalty} = \log \text{prior} \quad \text{MAP} \quad \left\| \begin{array}{l} \text{Maximum a} \\ \text{posteriori} \end{array} \right.$$
$$= -\frac{1}{2\sigma^2} \underline{w}^T \underline{w}$$

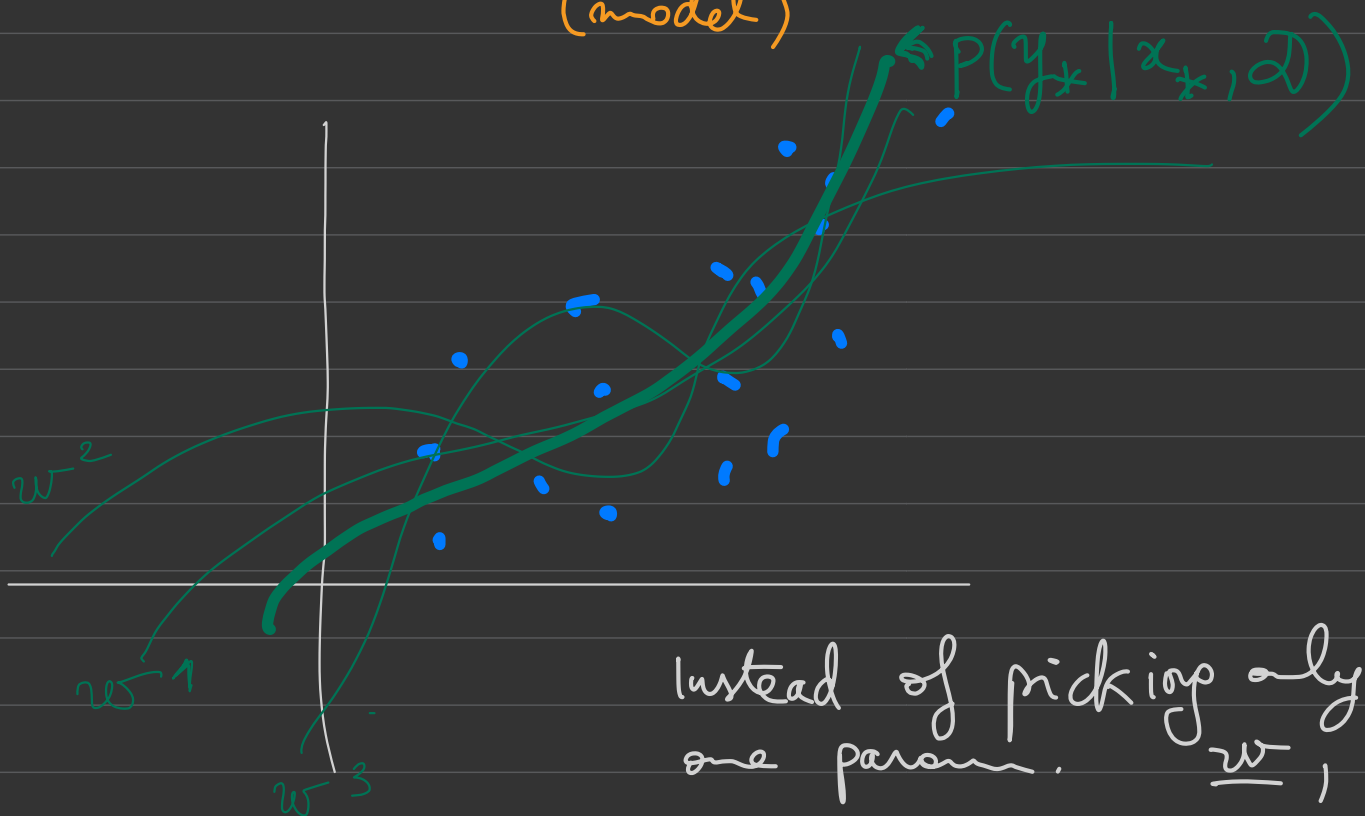
$L^2$  penalty, aka weight decay in ML/DL literature.

MAP is not quite Bayes because it's optim<sup>o</sup>  
Instead: we want a full predictive dist.

$$P(y_* | x_*, \mathcal{D}) = \int P(y_*, \underline{w} | x_*, \mathcal{D}) d\underline{w}$$

$$= \int P(y_* | \underline{w}, x_*, \cancel{\mathcal{D}}) P(\underline{w} | \cancel{\mathcal{D}}) d\underline{w}$$

$$= \int \underbrace{P(y_* | \underline{w}, x_*)}_{\text{likelihood (model)}} \underbrace{P(\underline{w} | \mathcal{D})}_{\text{posterior}} d\underline{w}$$



Instead of picking only  
 one param.  $\underline{w}$ ,

Bayesian model averaging takes them  
 all and weight them w.r.t. posterior.

Quick detour:

1) Likelihood of flipping  $n$  coins w/ prob.  $\lambda$  of tail.  $y_i = 1$  when tail

$$p(y_1, \dots, y_n | \lambda) = \prod_{i=1}^n p(y_i | \lambda) \\ = \lambda^{\sum y_i} (1-\lambda)^{n - \sum y_i}$$

Binomial ( $m | n, \lambda$ )

$$p(m | n, \lambda) = \binom{n}{m} \lambda^m (1-\lambda)^{n-m}$$

2) MLE for  $\lambda$ ?

$$\hat{\lambda}^{\text{MLE}} = \operatorname{argmax}_{\lambda} p(m | n, \lambda) = \frac{m}{n} = \bar{y}$$

3) Suppose you observed one tail.

What's the probab. that the next flip is tail again? (using MLE).

$$m = 1, n = 1 \Rightarrow \hat{\lambda}^{\text{MLE}} = \frac{1}{1} = 1$$



Bayes. approach

like  $p(m|m, \lambda)$

prior  $p(\lambda) = \text{Beta}(\lambda; a, b) \quad a, b > 0$

$\Rightarrow$  conjugate, so

post.  $p(\lambda | m, n) = \text{Beta}(\lambda; a+m, b+n-m)$

$$\hat{\lambda}^{\text{Bayes}} = E[\lambda | m, n] = \frac{a+m}{a+b+n}$$

$$a = b \quad \hat{\lambda}^{\text{Bayes}} < \lambda$$

predictive distr.:

$$p(y_* | x_*, \mathcal{D}) = \int p(y_* | \underline{w}, x_*) p(\underline{w} | \mathcal{D}) d\underline{w}$$

$\nwarrow$   
MLE chooses only one  
specif  $\underline{w}^0$ , i.e. chooses  
an approx. post  $\approx \delta(\underline{w}^0)$ .