

## BEST VALUES FOR MODEL PARAMETERS

### (15) Maximum Likelihood Estimation

Parameter Estimation: We often assume a parametric model in ml or stats, ie data is generated from some distribution  $P(x; \theta)$  or  $P(y|x; \theta)$  governed by params  $\theta$

Goal: Estimate  $\theta$  from observed data.

$\theta$ , that makes the observed data most likely

Likelihood function:  $l(\theta) = P(x_1, x_2, \dots, x_n | \theta)$

If samples are independent:  $L(\theta) = \prod_{i=1}^n P(x_i | \theta)$

MLE choose,  $\hat{\theta}_{MLE} = \arg \max_{\theta} l(\theta)$

To simplify math, we take log.  $\hat{\theta}_{MLE} = \arg \max_{\theta} l(\theta)$

Why?  
Why?

- Every ML model has some parameters associated with them.
- We don't know their values initially
- So we must estimate them from training data.

In logistic regression  $\rightarrow$  MLE: cross-entropy loss

In linear regression  $\rightarrow$  MLE: MSE

MLE = pure data fitting (no prior assumptions)

- \* Depending on what probability distribution you assume for data you get different loss functions
- \* Every loss function in ML and DL is just the negative log likelihood of a suitable probabilistic model.

Example: Linear Regression

1. Assume:  $y_i = w^T x_i + \epsilon_i$ ,  $\epsilon_i \sim N(0, \sigma^2)$  {Normal?}

2. Likelihood:

$$P(y_i | x_i, w) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - w^T x_i)^2}{2\sigma^2}\right)$$

3. Log likelihood:

$$\ell(w) = -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w^T x_i)^2 + c$$

4. Negative log likelihood:

$$J(w) = \frac{1}{2} \sum_{i=1}^n (y_i - w^T x_i)^2$$

This is the frickin MSE!

MLE work great when you've

- Plenty of data
- No noise
- Parameters are identifiable

MLE tends to overfit data (training) because it only tries to maximize fit - it doesn't penalize large or extreme parameter values.

fix?

⑯ Maximum A Posteriori Estimation. (MAP)

Now we take a Bayesian, constrainted approach

Instead of just maximising  $P(D|O)$  we also consider how likely the parameter itself is - using a prior distribution  $P(\theta)$

$$\text{Bayes Rule: } P(\theta | D) = \frac{P(D|\theta) P(\theta)}{P(D)}$$

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

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

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

$\leftarrow$  data fitting term       $\downarrow$  regularization term  
 (likelihood)

$$\left\{ \begin{array}{l} \text{MAP} = \text{MLE} + \text{Regularization} \end{array} \right\}$$

Gaussian Prior  $\rightarrow L_2$ , MAP works well for small data.

Laplace Prior  $\rightarrow L_1$ , Uncertain params.

- MLE: defines the data fit part of the loss.  
MLE gives us the core cost/loss function of the model.  
 $\rightarrow$  Best  $\theta$  maximising likelihood
- MAP: adds a prior on parameters (regularization)  
MAP gives loss + regularization term.  
 $\rightarrow$  Best  $\theta$  maximising posterior.

⑯ Bayesian Thinking: is a way of doing inference -  
ie, learning about unknown quantities based on data  
using Bayes Theorem.

$$\left\{ \text{Posterior} = \text{Prior} \times \text{Likelihood} \right\}$$

- Frequentist (MLE) : Fixed, but unknown constraints.
  - Single best  $\hat{\theta}$  (max likelihood)
  - It outputs a point estimate
- Bayesian (MAP) : random variables with probability distribution
  - Full distribution
  - Uses a prior  $\rightarrow$  output: updates belief.

{ Use the data - vs - update myself If I'm wrong }

## ⑦ Bayesian Estimation

Instead of finding a single best parameter value (like MAP / MLE)  
 we compute the entire probability distribution of the  
 parameters after seeing data.

So instead of a single  $\hat{\theta}$ , we have a full posterior dist

$P(\theta | D) \rightarrow$  Bayesian estimation.

Bayesian Linear Regression, model:  $y = w^T x + \epsilon$ ,  $\epsilon \sim N(0, \sigma^2)$

1. Starting with prior:  $w \sim N(0, T^2 I)$

2. Compute likelihood:  $P(D|w) \propto \exp\left(-\frac{1}{2\sigma^2} \|y - Xw\|^2\right)$

3. Posterior:  $P(w|D) \propto \exp\left(-\frac{1}{2\sigma^2} \|y - Xw\|^2\right) \exp\left(-\frac{1}{2T^2} \|w\|^2\right)$   
 by Bayes Theorem

This posterior is also gaussian.

4. We take log and negate:  $\text{loss}(w) = \|y - Xw\|^2 + \lambda \|w\|^2$

≈ Ridge Regression - the MAP estimate of the Bayesian Model.

In easy terms →

- Prior: What you believe before data
- Likelihood: How data pushes your belief
- Posterior: Your updated belief
- MAP: Most likely belief, after update

Bayesian Estimation: Full range of beliefs after update.

⇒ Why use Bayesian Estimation

- Uncertainty estimation
- Small data robustness
- Regularization built-in
- Sequential Updating

## ⑧ Complete Distributions

A probability distribution describes how likely different outcomes are. Tells how randomness behaves.

→ Discrete distributions: events that take specific values

① Bernoulli distribution: Used for binary outcomes

$$P(X=x) = p^x (1-p)^{1-x}, \quad x \in \{0,1\}$$

→ logistic regression: Binary classification

② Binomial Distribution: Sum of n Bernoulli trials

$$P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}$$

→ Continuous distribution: events that take any real value

③ Normal / Gaussian distribution: classic bell curve

$$P(x) = 1/\sqrt{2\pi\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$