

$(x_i, y_i) \sim$  dist of data iid

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

model:  $x \mapsto$  distribution on  $y$  params  $\theta$

$\max_{\theta} P(D|\theta)$  vs.  $\max_{\theta} P(\theta|D)$

MLE

MAP

frequentist

max a posteriori estimation

prior  $P(\theta)$

collect  $D$

update prior  $P(\theta|D)$

Bayesian

Bayes:

$$\log P(\theta|D) = \log P(D|\theta) + \log P(\theta) - \log P(D)$$

If  $P(\theta) \equiv \text{const} \Rightarrow$  equivalent MLE  $\Leftrightarrow$  MAP

uninformative prior

Does there exist a prior  $P(\theta)$  uninformative over  $\mathbb{R}^d$ ?

Improper prior

MLE training of a NN - MAP Estimation  
 w/ noninformative improper prior  
 Weight Decay

$$\min_{\theta} \underbrace{-\log P(D|\theta)}_{\mathcal{L}_{CE}(D|\theta)} + \|\theta\|_2^2$$

Bayesian perspective: prior  $\log P(\theta) = -\|\theta\|_2^2$   
 $P(\theta) \sim \mathcal{N}(0, I)$

Dist  $P(x|\theta)$   $P_{\theta}(x)$

How sensitive is it to changes in  $\theta$ ?

$$\nabla_{\theta} \log P(x|\theta) = 0 \text{ at a sdn to training}$$

Instead,  
 look at  $D_{\theta}^2 \log P(x|\theta)$

$$\mathbb{E}_{x \sim P_{\theta}} D_{\theta}^2 \log P(x|\theta) = \mathbb{E}_x \underbrace{\nabla \log P(x|\theta) \nabla \log P(x|\theta)^T}_{\text{Fisher Information of } P_{\theta}}$$

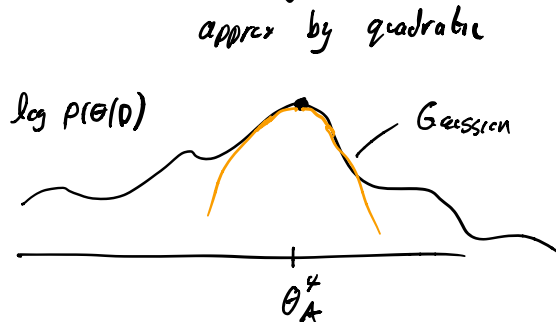
large diagonal entries have large information content

Fisher info is covariance<sup>-1</sup> of  $\nabla \log P$

$$\log P(\theta|D) = \log P(D_B|\theta) + \log P(\theta|D_A) \dots$$

Laplace approximation  
 of a posterior

$\log P(\theta)$   
 \_\_\_\_\_  
 prior



$$\log P(\theta | D_n) = \log P(\theta_n^* | D_n) + \frac{1}{2}(\theta - \theta_n^*)^t H (\theta - \theta_n^*)$$

$$H = \mathbb{E} \left[ \nabla_{\theta} \log P \nabla_{\theta} \log P^t \right]$$

$$\mathbb{E}_{\mathcal{X}} \left[ (\nabla_{\theta} \log P(y_i | x_i, \theta))^2 \right] = \text{diag}(F)$$

$$F = \mathbb{E} \left[ \nabla_{\theta} \log P(x_i | \theta) \nabla_{\theta} \log P(x_i | \theta)^t \right]$$

$$(x_i, y_i) \quad P(x_i, y_i | \theta) = P(y_i | x_i, \theta) P(x_i)$$