

# Week 11 events

No MLPR Lectures

(last lecture Thurs 21 Nov.)

Ed - Intelligence have two  
(Links also in "week 11" of mlpr notes) events!

## Mini NeurIPS

6pm Wed 27 Nov, AT LT 2

[tinyurl.com/mini-neurips-2019](http://tinyurl.com/mini-neurips-2019)

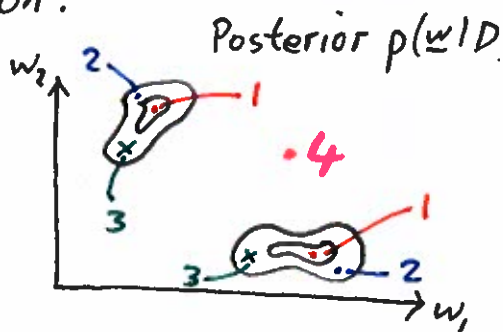
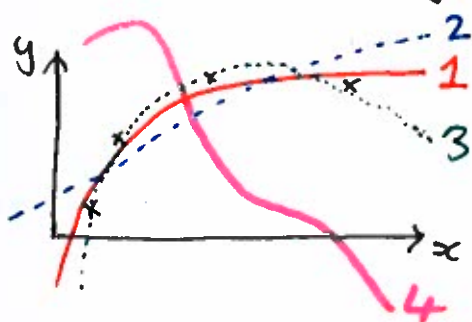
## Biases, failure + fairness in AI

6pm Fri 29 Nov, AT LT 5

[to-err-is-machine.eventbrite.co.uk](https://to-err-is-machine.eventbrite.co.uk)

# Approximate Bayesian Inference

E.g. neural net regression:

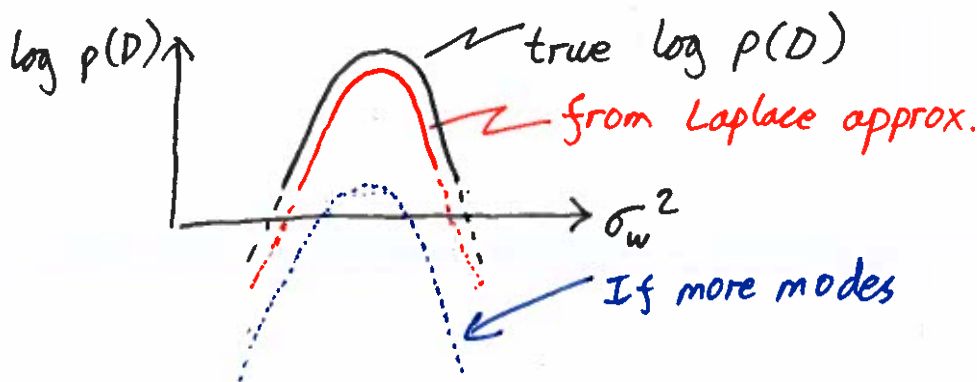


Laplace:  $p(w|D) \approx N(w; w^*, H^{-1})$

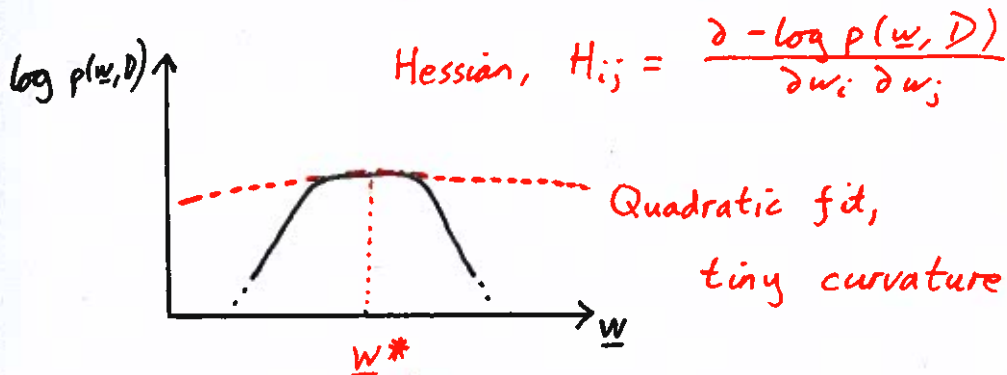
Prediction: fitting one mode might be ok.

If 2 Gaussian modes,  $p(D)$  half true answer

Setting hyper-parameters might work:



# Laplace Approximation



Posterior

Found with optimizer

$$p(\underline{w} | D) \approx N(\underline{w}; \underline{w}^*, H^{-1})$$

Marginal Likelihood

$$p(D) = \frac{p(\underline{w}^*, D)}{p(\underline{w}^* | D)} \approx \frac{p(\underline{w}^*, D)}{N(\underline{w}^*; \underline{w}^*, H^{-1})}$$



## Quiz

Is  $p(D)$  approx.

- A) Too big
- B) Too small
- C)  $\approx$  Correct
- Z) ???

L28 (3)

2019 L27 (8)

Today other approximations  
to  $p(\underline{w} | D)$

1) Monte Carlo

"random sampling"

2) "Variational methods"

make Bayesian inference optimization

Reminder Bayes' Rule

$$p(\underline{w} | D) = \frac{p(D | \underline{w}) p(\underline{w})}{p(D)}$$

$\int p(D | \underline{w}) p(\underline{w}) d\underline{w}$

# Monte Carlo

$$P(y=1 | \underline{x}, \mathcal{D}) = \mathbb{E}_{p(\underline{w} | \mathcal{D})} [\sigma(\underline{w}^T \underline{x})]$$

$$\approx \frac{1}{S} \sum_s \sigma(\underline{w}^{(s)T} \underline{x}), \quad \underline{w}^{(s)} \sim p(\underline{w} | \mathcal{D})$$

Approximately using "MCMC" How?  
(not this course)

# Importance Sampling

$$\int g(x) p(x) dx = \int \left[ g(x) \frac{p(x)}{q(x)} \right] q(x) dx$$

$\nearrow q(x) \neq 0$   
when  $p(x) \neq 0$

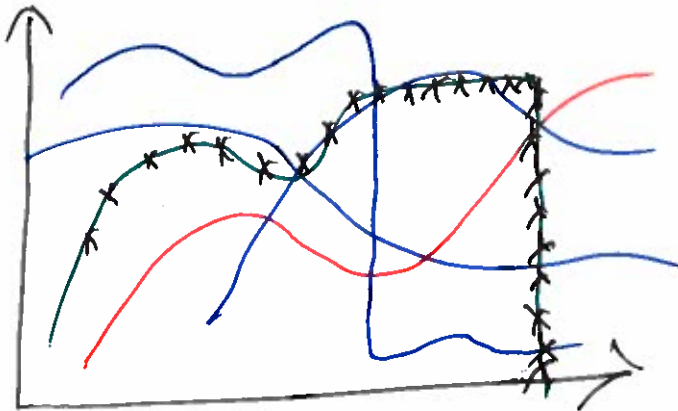
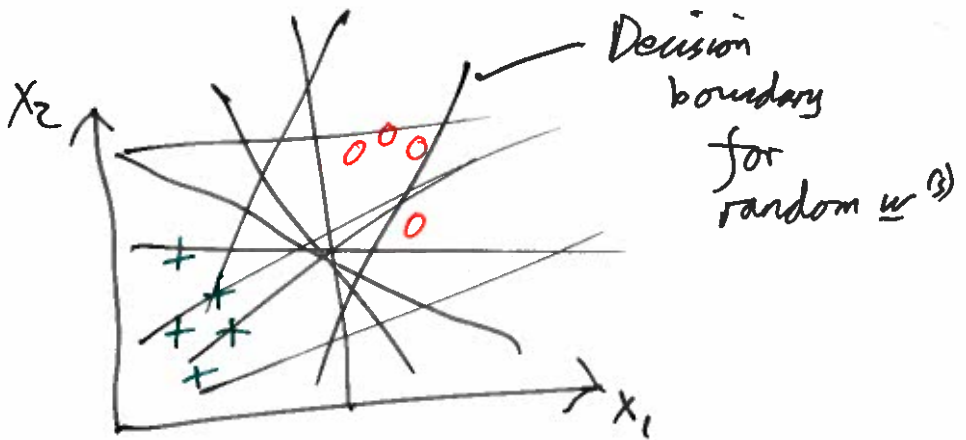
$$\mathbb{E}_p [g(x)] = \mathbb{E}_q \left[ g(x) \frac{p(x)}{q(x)} \right] \approx \frac{1}{S} \sum_s g(x) \frac{p(x)}{q(x)}$$

$x \sim q$

For logistic regression (small test cases)

$$\underline{w}^{(s)} \sim q(\underline{w}) = \text{prior } p(\underline{w})$$

$$P(y=1 | \underline{x}, \mathcal{D}) \approx \frac{\frac{1}{S} \sum_{s=1}^S \sigma(\underline{w}^{(s)T} \underline{x}) \frac{P(\mathcal{D} | \underline{w}^{(s)}) p(\underline{w}^{(s)})}{p(\underline{w}^{(s)})}}{\frac{1}{S} \sum_{s=1}^S P(\mathcal{D} | \underline{w}^{(s)}) \frac{p(\underline{w}^{(s)})}{p(\underline{w}^{(s)})}} \frac{1}{P(\mathcal{D})}$$



## Variational Method

Another way to fit approximation to  $p(\underline{w} | D)$

$$p(\underline{w} | D) \approx q(\underline{w}; \alpha)$$

For us  $q(\underline{w}; \alpha) = \mathcal{N}(\underline{w}; \underline{m}, V)$

"Variational parameters"

Have optimization problem

Fit  $\alpha$ , need cost function

Measure difference between  $p(\underline{w} | D)$  and  $q(\underline{w})$

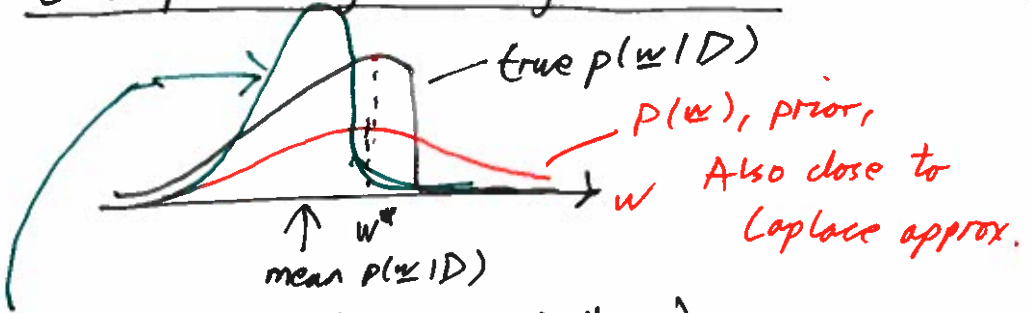
Often Kullback-Leibler Divergence (KL)

$$D_{KL}(r \parallel s) = \int r(\underline{z}) \log \frac{r(\underline{z})}{s(\underline{z})} d\underline{z}$$

$$\geq 0 \quad (\text{Gibbs' inequality})$$

It isn't a distance:  $D_{KL}(r \parallel s) \neq D_{KL}(s \parallel r)$

# Example Logistic Regression $N=1$



Minimize  $D_{KL}(p(w|D) || q)$

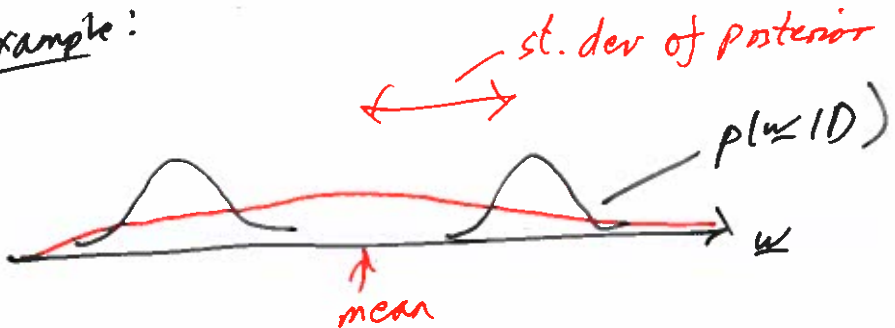
⇒ Match mean & covariance of posterior

Don't normally minimize  $D_{KL}(p(w|D) || q)$

1) It's harder (check you know why)

2) Often not a good idea:

Example:





Minimizing  $D_{KL}(q \parallel p(\underline{w}|D))$

$$D_{KL}(q \parallel p(\underline{w}|D))$$

$$= \int q(\underline{w}; \alpha) \log \frac{q(\underline{w}; \alpha)}{p(\underline{w}|D)} d\underline{w}$$

$$= \underbrace{-\int q(\underline{w}; \alpha) \log p(\underline{w}|D) d\underline{w}}_{\text{good}} + \underbrace{\int q(\underline{w}; \alpha) \log q(\underline{w}; \alpha) d\underline{w}}_{\text{Entropy of } q}$$

good  $q(\underline{w}; \alpha)$  is big  
when  $\log p(\underline{w}|D)$  is big

- Entropy of  $q$   
"spread out"

Really bad if  $q(\underline{w}; \alpha)$  is big  
when  $p(\underline{w}|D)$  is tiny