$$\boxed{\text{TRAINING}}$$

$\underline{x}$ features  $\qquad$  $D$ data

$\underline{y}$ labels

$\underline{w}$ weights

$$\underline{NN = model \quad P(\underline{y} \mid \underline{x}, \underline{w}, D)}$$

"STANDARD" NN : learn ML estimator

$$\underline{w}^* = \arg\max_w \log \mathcal{L}(D \mid \underline{w})$$

$$= \arg\max_w \sum P(y_i \mid x_i, w_i, D)$$

$\uparrow$

NO GLOBAL MINIMUM $\forall$ $\boxed{D}$ !

⚑ or MAP w. prior

$$w^{MAP} = \arg\max_w \log \mathcal{L}(D \mid \underline{w}) + \log p(\underline{w})$$

p(7) ~ 56%     p(2) ~ 43%



p(7)

1

p(2)

0.5

DATA

SIGMOID INPUT

DATA

$x$

$W^b$

POINT - LIKE

$y$

DISTRIBUTIONS

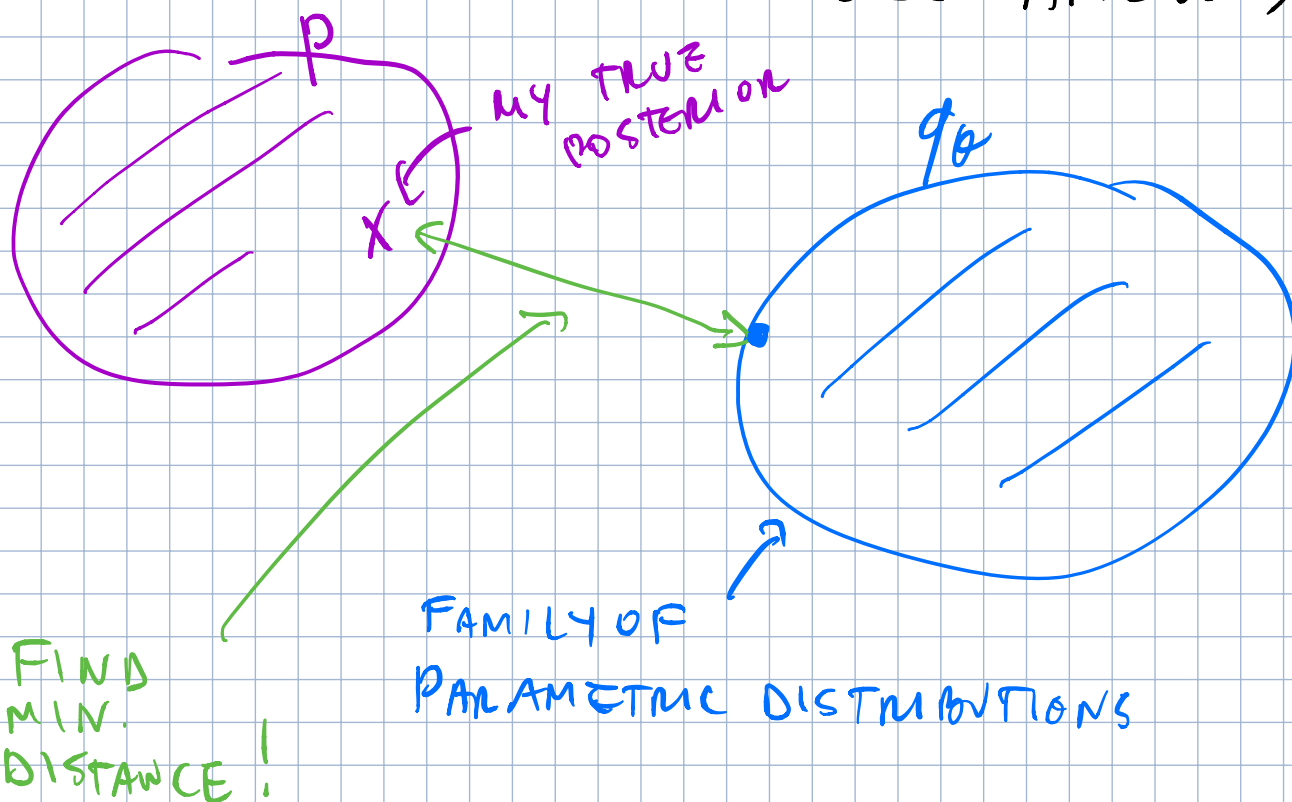NOW : LEARN $\boxed{p(\underline{w} \mid \mathcal{D})}$
POSTERIOR ON WEIGHTS

$$p(\underline{y} \mid \underline{x}) = \mathbb{E}_{p(\underline{w} \mid \mathcal{D})}\left[ p\left( \underline{y} \mid \underline{x}, \underline{w} \right) \right]$$

$p(\underline{w}, \mathcal{D})$ full unknown distribution
(COULD DO HMC ...)



p

MY TRUE POSTERIOR

$x^c$

$q_\theta$

FIND MIN. DISTANCE!

FAMILY OF PARAMETRIC DISTRIBUTIONS

$q(\underline{w} \mid \theta)$   PARAMETRIC DISTRIBUTION

THAT "BEST APPROXIMATES" $p(\underline{w} \mid \mathcal{D})$

# KL DIVERGENCE

$$KL\left[f \| g\right] = \int f(x) \log \frac{f(x)}{g(x)} dx$$

MINIMIZE:

$$KL\left[q(\underline{w}|\theta) \| p(\underline{w}|D)\right] =$$

$$= \int dw \; q(w|\underline{\theta}) \log \frac{q(\underline{w}|\underline{\theta})}{p(\underline{w}|D)} \quad *$$

BAYES: $\quad p(\underline{w}|D) \propto \mathcal{L}(D|w) \, p(w)$

$$= \int dw \; q(\underline{w}|\underline{\theta}) \log \frac{q(\underline{w}|\underline{\theta})}{\mathcal{L}(D,\underline{w}) \, p(\underline{w})}$$

$$= \int dw \; q(\underline{w}|\underline{\theta}) \log \frac{q(\underline{w}|\underline{\theta})}{p(\underline{w})} - \int dw \; q(w,\theta) \times \log \mathcal{L}(D|w)$$

$$= KL\left[q \| p(\underline{w})\right] + \mathbb{E}_q\left[-\log \mathcal{L}\right]$$

OPTIMIZE:

$$F(D, \theta) = KL\left[ q_\theta(\underline{w} | D) \| P(\underline{w}) \right]$$

$$+ \mathbb{E}_{q_\theta}\left[ -\log \mathcal{L}(D | \underline{w}) \right]$$

$$\approx \sum_{(i) \in \text{SAMPLES}} \log q_\theta(\underline{w}^{(i)} | D) - \log p(\underline{w}^{(i)}) - \log \mathcal{L}(D | \underline{w}^{(i)})$$

$$\underbrace{\phantom{\log q_\theta(\underline{w}^{(i)} | D)}}_{\text{TARGET}} \quad \underbrace{\phantom{\log p(\underline{w}^{(i)})}}_{\text{PRIOR}} \quad \underbrace{\phantom{\log \mathcal{L}(D | \underline{w}^{(i)})}}_{\text{LOSS}}$$

$$= -\text{ELBO} \quad (\text{Evidence Lower Bound})$$

- BERNOULLI VIA DROPOUT     1506.02142

Proof that $\exists$   $\underline{p(\underline{w})}$ such that

$$F \simeq \underbrace{\sum \mathcal{L}_{\tilde{\imath}}(D, w)}_{\text{DATA}} - \ell \underbrace{\sum \| w_i \|^2}_{\text{WEIGHTS}}$$

$\quad$ + DROPOUT AFTER EACH LAYER

NOT SO "BAYESIAN" ...

- BAYES BY BACKPROP     1505.05424

Sample weights in back-propagation
(Once for batch)

$$F \simeq \sum \log q_\theta(w^{(i)}|D) - \log p(w^i) - \log \mathcal{L}(D | w^{(i)})$$

↑
<u>Millions of samples needed</u> + same sample for
$\qquad\qquad\qquad\qquad\qquad\qquad$ every batch

- LOCAL REPARAMETRIZATION TRICK     1506.02557

Sample layer activations instead of $w^{(i)}$

- FLIPOUT    1803.04386

  Samples weight independently within a mini-batch.

- NORMALIZING FLOWS    1802.04908

  Differentiable mapping to complex distributions

---

SEE NOTEBOOK!

tfp:    FLIPOUT + LOCAL R.T.

# PREDICTIONS & UNCERTAINTY

$$P(\underline{y}^*, \underline{x}^*) = \mathbb{E}_{q_\theta}\left[\underline{P}(\underline{y}^*, \underline{x}^*)\right]$$

$$\simeq \frac{1}{N_{samples}} \sum P\left(\underline{y}^* \mid \underline{x}^*, \underline{w}^{(i)}\right)$$

Each time draw 1 sample

$$Var_q\left(P(\underline{y}^*, \underline{x}^*)\right) = \mathbb{E}\left[\underline{y}^*, \underline{y}^{*T}\right] - \mathbb{E}_q(\underline{y}^*)\,\mathbb{E}_q(y^*)^T$$

~ Train last layer to learn $\mu, \sigma \ \forall$ label

~ Then:

$$Var_q\left(P(\underline{y}^*, \underline{x}^*)\right) \simeq \frac{1}{N}\sum_i \sigma_i^2 + \frac{1}{N}\sum (\mu_i - \bar{\mu}_i)^2$$

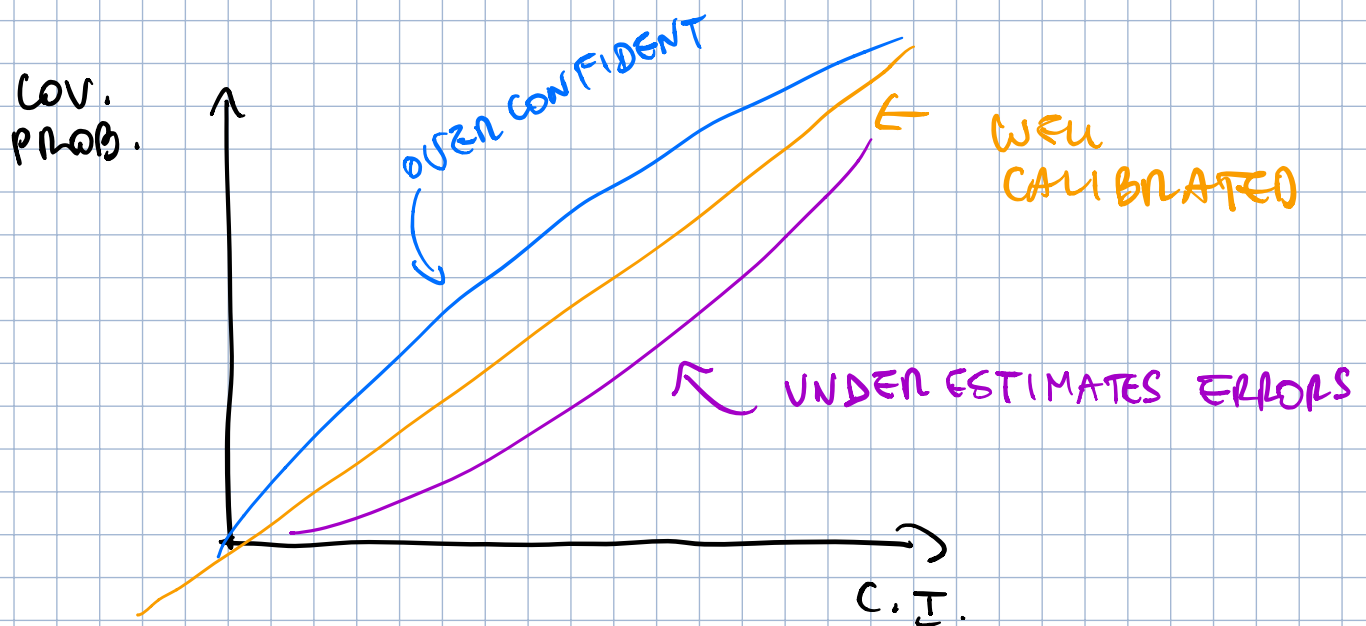Aleatoric                          Epistemic

?

## CALIBRATION

1708.08843
1911.08508

- Get an uncertainty

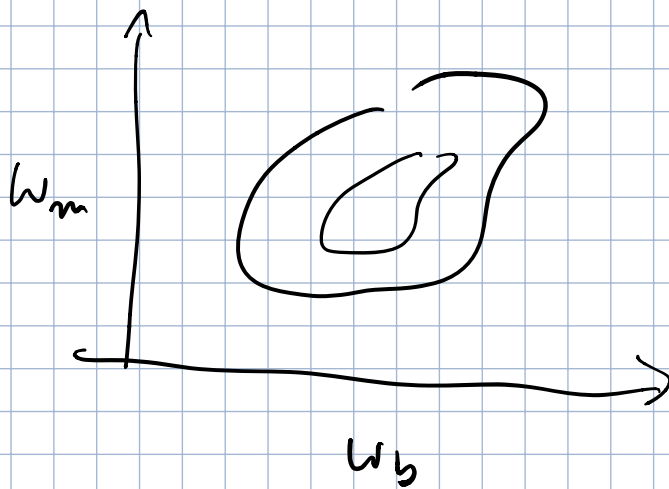- Coverage probability : fraction of test examples where the true value lies within a confidence interval



- Tune parameters to get correct result (e.g. Dropout rate)

- Methods for calibrating after training

Say now the last layer gives you some parameters $(\Omega_b, \Omega_m, A_s \qquad )$

Can draw samples :



$\omega_m$

$\omega_b$

— CONDITIONED ON DATA

— WHERE IS PRIOR?

— NOT "EXACT" POSTERIOR BUT $q_\theta$