

DD2412 - Deep Learning Advanced Course

Labelled data $D = \{(\bar{x}_i, y_i)\} \quad i \in 1:n$

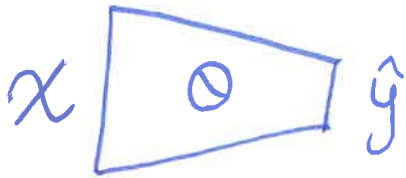
\uparrow Input \uparrow Label

Classification $y \in 1, 2, 3 \dots k$

Regression $y \in \mathbb{R}$

Weights $\theta \in \Theta$

$$f_{\theta}: x \mapsto \hat{y}$$



$$\theta^* = \operatorname{argmin}_{\theta} L(D)$$

Find θ where the loss function over all D is minimal

$$L(D) = \sum_{i=1}^n \ell(f_{\theta}(x_i), y_i) + \Omega(\theta)$$

\downarrow label

\uparrow Regularization

eg: $\ell = \text{MSE}$
 $\Omega = \|\theta\|^2 (L_2)$

Probabilistic discriminative learning

$$\theta^* = \operatorname{argmax}_{\theta} P(\theta | D)$$

What parameters are most likely to match D

$$\Rightarrow \operatorname{argmax}_{\theta} \frac{P(D|\theta)P(\theta)}{P(D)}$$

$$\operatorname{argmax}_{\theta} \prod_i P(x_i, y_i | \theta) P(\theta)$$

$$\operatorname{argmax}_{\theta} \sum \log(P(x_i, y_i | \theta)) + \log(P(\theta))$$

$\underbrace{\hspace{10em}}_{P(y_i | x_i, \theta)}$

Bayesian Modeling

$$P(\theta | D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

\uparrow Posterior \uparrow Likelihood \uparrow evidence \uparrow Prior

$$\theta^* = \operatorname{argmax}_{\theta} P(D|\theta)P(\theta)$$

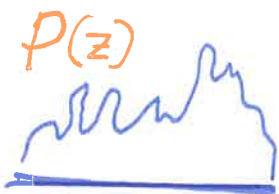
Usually the evidence is hard to calculate so it needs to be approximated

Can be used to make new likelihood predictions of new datapoints

Variational Inference

General Steps

- we want to find a distribution P (eg $P(\theta|D)$)
- but the distribution is hard to find directly
- Settle with some approximation of P from a simpler family of distributions, Q , parameterized by ω
- find a parameter ω^* that minimizes difference between Q and P (KL divergence)



find a ω^* that is as similar to P as possible

↓
Lowest
KL-divergence

KL-divergence

$D_{KL}(P||D)$ or $KL(P||D)$ measure to compare P to D

$$KL(P||D) = 0 \Leftrightarrow P = D$$

$$\begin{aligned} KL(P||D) &= E_P [I_Q(x) - I_P(x)] \\ &= \sum (-\log Q(x) + \log P(x)) P(x) \\ &= \sum P(x) \log \frac{P(x)}{Q(x)} \\ &= H(P, Q) - H(P) \end{aligned}$$

(moment)

$KL(P||D) \Rightarrow$ Mean matching

$KL(D||P) \Rightarrow$ Mode matching

↳ usually the one used

Reminders:

Information $I_P(x) = -\log P(x)$

Entropy $H(P) = E_P [I_P(x)] = -\sum \log(P) P(x)$

↳ 0 \Rightarrow 100% predictable

↳ max \Rightarrow uniform P

three observations.

KL is not symmetric.

for some $z \in P$

- $P(z)$ and $Q(z)$ is high \Rightarrow KL low
- $P(z)$ is low $Q(z)$ is high \Rightarrow KL still low
- $P(z)$ is high $Q(z)$ low \Rightarrow KL high

$\Rightarrow Q$ must be high whenever P is
but not necessarily other way around

Steps of Variational Inference

given observations X and hidden variables Z and we are interested in the true posterior, $P(Z|X)$

Use bayes rule $P(Z|X) = \frac{P(X|Z)P(Z)}{P(X)}$

Consider a simpler distribution $Q_\omega(Z)$, match Q to $P(Z|X)$

$\hookrightarrow P(X)$ is hard to find

$$P(Z|X): \min_{\omega} D_{KL}(Q_{\omega}(Z) \| P(Z|X))$$

\hookrightarrow optimization instead of integration.
 $\hookrightarrow Q$ is differentiable

Discriminative Bayesian modeling using VA

1. have training data $D = \{(x_1, y), \dots, (x_n, y_n)\}$

2. Assume prior distribution $P(\theta)$

3. Design likelihoods $P(y|x, \theta)$

4. Decide on approximation distribution $Q_{\omega}(\theta)$

5. Find ω with the following objective

$$\max_{\omega} \sum (\underbrace{Q_{\omega}(\theta)}_{\sum_i} \log P(y_i | x_i, \theta)) - D_{KL}(Q_{\omega}(\theta) \| P(\theta))$$

Predictive Uncertainty.

We want the model to give out a certainty with its prediction both classification & regression

Data points that the model is uncertain about are the ones where it should focus on more

Why?

Small dataset, noisy in/out-put, incomplete input, covariate shift

Epistemic Vs Aleatoric

↳ Lack of Knowledge

- not enough data
- Networks too simple
- Bad optimizer
- Not same distribution

↳ Model Uncertainty

↳ Distributional uncertainty

↳ Unobtainable Knowledge

- Label noise
- Measurement precision
- measurement noise
- Class definition overlap

↳ Homoscedastic Uncertainty

- Absolute uncertainty
- constant uncertainty

↳ Heteroscedastic Uncertainty

- Relative uncertainty
- input dependent

How to model predictive uncertainty?

How to model Aleatoric Uncertainty?

We don't know the noise in the data

$$y = y + \text{noise}$$

↳ ϵ independent of x ? EVS $E(x)$

$$\mathcal{N}(0, \sigma^2)$$
$$\mathcal{N}(0, \sigma(x)^2)$$

Assume a distribution (typically gaussian)

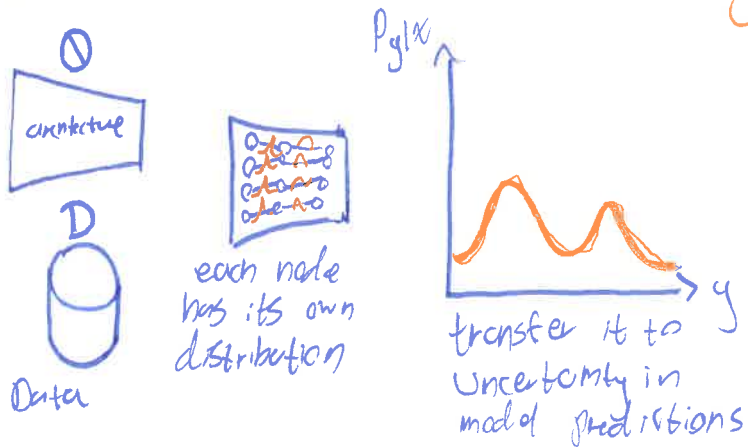
1 have the model f_θ output a distribution over y

2 find θ with MLE

3 have loss according to

$$l(f_\theta(x_i), y_i) = \frac{\log \sigma^2}{2} + \frac{(f_\theta(x_i) - y_i)^2}{2\sigma^2} + C$$

How to model epistemic uncertainty



Frequentist approach

train model multiple times and model the amount of disagreement between models.

Bayesian approach

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

get $P(y|x, D)$

for discriminative

As in MAP

assume a prior $P(\theta)$

devise a likelihood function $p(y|x, \theta)$

assuming D is iid

full posterior over θ $P(\theta|D) = \frac{\prod_i P(y_i|x_i; \theta) P(\theta)}{P(D)}$

Variational Inference for Deep uncertainty Estimation

↳ Example

1 have data D

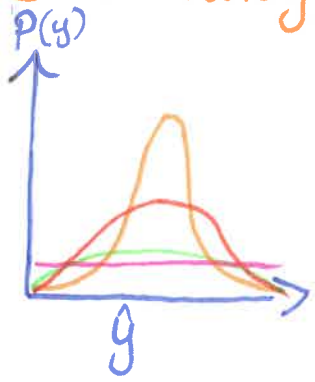
2 have prior $P(\theta) = \mathcal{N}(\mu, I)$

3 Design likelihood $P(y_i|x_i, \theta) = f_{NN}(x_i; \theta)$ - softmax

4 assign approximation function $Q_{\omega}(\theta) = \mathcal{N}(\mu, \sigma^2 I)$

5 find ω that maximises $*$

Uncertainty Estimation



Aleatoric Uncertainty → irreducible → Maximum Likelihood estimation of noise

Epistemic Uncertainty
Decreases with data

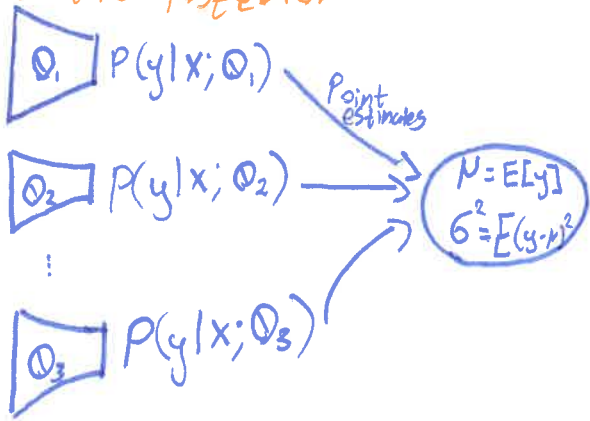
Bayesian Modeling
 $P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$

Variational approximation
 $\min_{\omega} D_{KL}(Q_{\omega}(z) || P(z|\nu))$

Ensemble Methods

Train S independent networks with bootstrapped data.

Have each network act as samples of the posterior



- + Not as many parameters
- + Model structure is unchanged
- + Parallelizable
- evaluation is inefficient
- lots of memory

SWAG

take snapshots of the model at different steps of the training and treat each estimate as its own network.

Distillation

train an ensemble of networks, use them to label data. Train a student network on the larger dataset.

Evaluating Uncertainty

Proper Scoring

Imagine a scoring function

$S(P_0(y|x), (x_i, y_i))$ that measures the quality of a predictive distribution with underlying data.

Define $S(P_0, P) = \iint P(x, y) S(P_0(x, y)) dx dy$
 S is the proper scoring if

$$S(P_0, P) \leq S(P, P)$$

S is strictly proper iff

$$P_0(y|x) = P(y|x)$$

Calibration

Measure statistical consistency between predictive and observed distribution

Log likelihood

Log likelihood

$$S(P_0, (x_i, y_i)) = \log(P_0(y=y_i|x_i))$$

Brier score

$$S(P_0, (x_i, s_i)) = \text{MSE}(P_0(y|x_i), s_i)$$

Out of Distribution Detection

What should the model do when test and training data come from different distributions. Possibility is to abstain from prediction.

