

Probabilistic Reasoning in Deep Learning

Dr Konstantina Palla, PhD
palla@stats.ox.ac.uk



September 2017

Deep Learning Indaba, Johannesburg

OVERVIEW OF THE TALK

- Basics of Bayesian Inference

OVERVIEW OF THE TALK

- Basics of Bayesian Inference
- Bayesian reasoning inside DNNs.

OVERVIEW OF THE TALK

- Basics of Bayesian Inference
- Bayesian reasoning inside DNNs.
 - > The **Bayesian paradigm** as an approach to account for uncertainty in the DNNs.

OVERVIEW OF THE TALK

- Basics of Bayesian Inference
- Bayesian reasoning inside DNNs.
 - > The **Bayesian paradigm** as an approach to account for uncertainty in the DNNs.
- Connection to other statistical models (**nonparametric** models).

OVERVIEW OF THE TALK

- Basics of Bayesian Inference
- Bayesian reasoning inside DNNs.
 - > The **Bayesian paradigm** as an approach to account for uncertainty in the DNNs.
- Connection to other statistical models (**nonparametric** models).
 - > How kernels and Gaussian Processes fit in the picture of NNs.

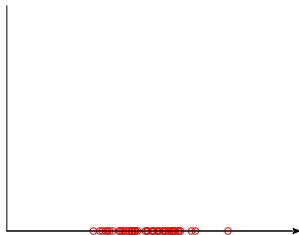
Bayesian Inference: Basics

BAYESIAN INFERENCE: BASICS

Problem: How was it **generated**?

Data: Observations

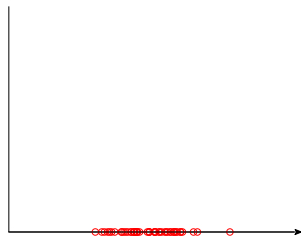
$$\mathbf{Y} = y^{(1)}, y^{(2)}, \dots, y^{(N)}$$



BAYESIAN INFERENCE: BASICS

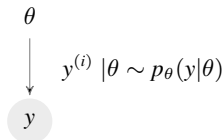
Data: Observations

$$\mathbf{Y} = y^{(1)}, y^{(2)}, \dots, y^{(N)}$$



Problem: How was it **generated**?

Answer: Make assumptions \rightarrow distributional



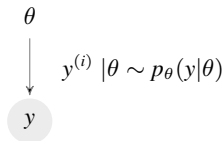
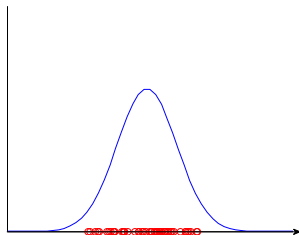
BAYESIAN INFERENCE: BASICS

Data: Observations

$$\mathbf{Y} = y^{(1)}, y^{(2)}, \dots, y^{(N)}$$

Problem: How was it **generated**?

Answer: Make assumptions \rightarrow distributional



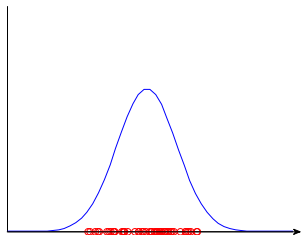
$$y^{(i)} | \theta \sim \mathcal{N}(y; \mu, \sigma^2)$$

$$\theta = \{\mu, \sigma^2\}$$

BAYESIAN INFERENCE: BASICS

Data: Observations

$$\mathbf{Y} = y^{(1)}, y^{(2)}, \dots, y^{(N)}$$

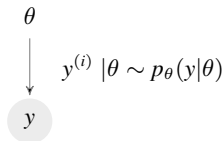


$$y^{(i)} | \boldsymbol{\theta} \sim \mathcal{N}(y; \mu, \sigma^2)$$

$$\boldsymbol{\theta} = \{\mu, \sigma^2\}$$

Problem: How was it **generated**?

Answer: Make assumptions \rightarrow distributional



likelihood: the probability of the data under the model parameters θ

$$\begin{aligned} p(\mathbf{Y}|\boldsymbol{\theta}) &= \prod_{i=1}^N p(y^{(i)}|\boldsymbol{\theta}) \\ &= \prod_{i=1}^N \mathcal{N}(y^{(i)}|\mu, \sigma^2) \end{aligned}$$

Bayesian Principle

All forms of uncertainty should be expressed by means of probability distributions.

Prior: $\theta \sim p(\theta)$

Our prior belief before seeing the data

Posterior: Our updated belief after having seen the data

θ



y

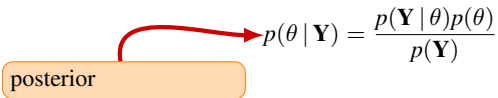
Bayes' Rule

$$p(\theta|\mathbf{Y}) = \frac{p(\mathbf{Y}|\theta)p(\theta)}{p(\mathbf{Y})}$$

Baye's rule

$$p(\theta | \mathbf{Y}) = \frac{p(\mathbf{Y} | \theta)p(\theta)}{p(\mathbf{Y})}$$

Baye's rule



posterior $p(\theta | \mathbf{Y}) = \frac{p(\mathbf{Y} | \theta)p(\theta)}{p(\mathbf{Y})}$

The posterior is proportional to **likelihood** \times **prior**

BAYESIAN INFERENCE: BASICS

Baye's rule

likelihood

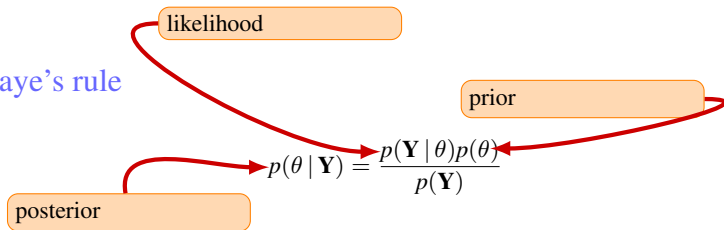
$$p(\theta | \mathbf{Y}) = \frac{p(\mathbf{Y} | \theta)p(\theta)}{p(\mathbf{Y})}$$

posterior

The posterior is proportional to **likelihood** \times **prior**

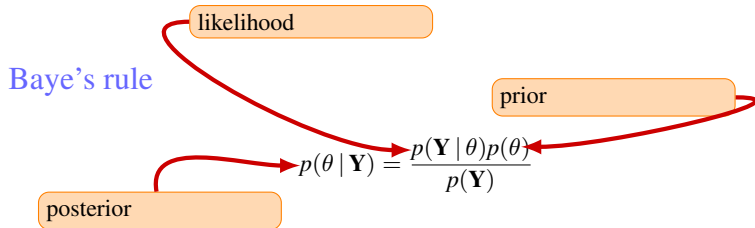
BAYESIAN INFERENCE: BASICS

Baye's rule



The posterior is proportional to **likelihood** \times **prior**

BAYESIAN INFERENCE: BASICS

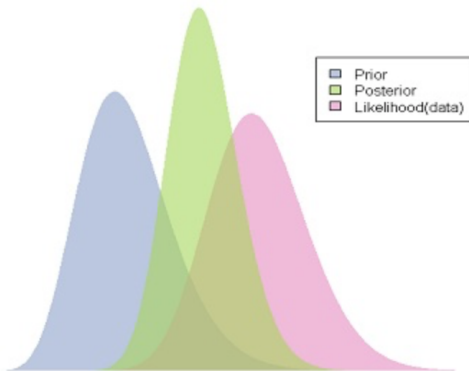


The posterior is proportional to **likelihood** \times **prior**

Update my **prior beliefs** after I see the **data**

BAYESIAN INFERENCE: BASICS

$$P(\theta|\mathbf{Y}) \propto P(\mathbf{Y}|\theta)P(\theta)$$



Make predictions

For unseen y^* marginalise!

$$P(y^* | \mathbf{Y}) = \int P(y^* | \theta) P(\theta | \mathbf{Y}) d\theta$$

Make predictions

For unseen y^* marginalise!

$$P(y^*|\mathbf{Y}) = \int P(y^*|\theta)P(\theta|\mathbf{Y})d\theta$$

The prediction does not depend on a point estimate of θ but it is expressed as a weighted average over all the possible values of θ

Why is this useful?

> Uncertainty is taken into account

No point estimates BUT distribution over the unknown θ ; $p(\theta|\mathbf{Y})$

Why is this useful?

> Uncertainty is taken into account

No point estimates BUT distribution over the unknown θ ; $p(\theta|\mathbf{Y})$

> Natural handling of missing data

- a probability distribution is estimated for each missing value



$$\theta \text{ ?}$$
$$\theta \sim p(\theta|\mathbf{Y})$$

Why is this useful?

> Uncertainty is taken into account

No point estimates BUT distribution over the unknown θ ; $p(\theta|\mathbf{Y})$

> Natural handling of missing data

- a probability distribution is estimated for each missing value



$$\theta \quad ?$$
$$\theta \sim p(\theta|\mathbf{Y})$$

> Addresses issues like regularization, overfitting \rightarrow **useful in Deep Neural Networks.**

Bayesian Reasoning in (Deep) Neural Networks

Bayesian Reasoning in (Deep) Neural Networks

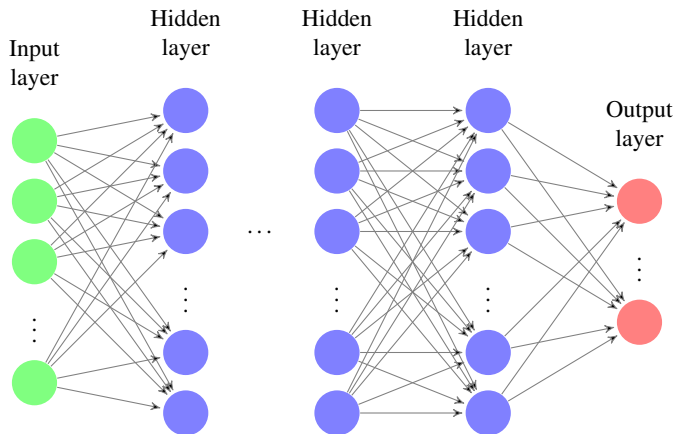
Overview of (D)NNs

Bayesian Reasoning in (Deep) Neural Networks

Overview of (D)NNs

Bayesian Reasoning in DNNs

OVERVIEW OF (DEEP) NEURAL NETWORKS

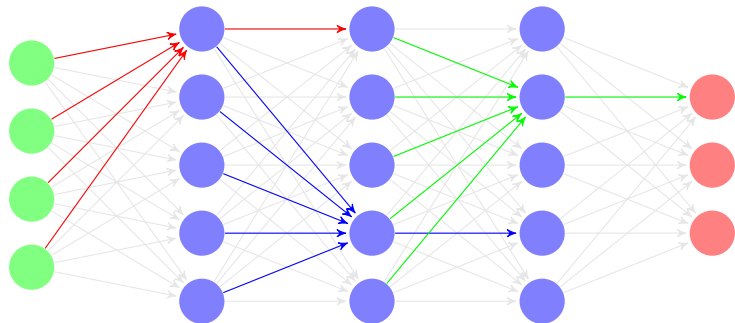


What is Deep Learning

A framework for constructing *flexible* models → a way of constructing outputs using functions on inputs.

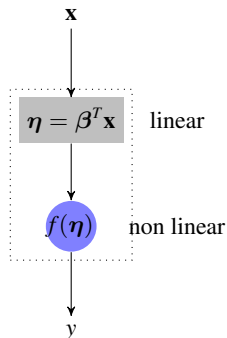
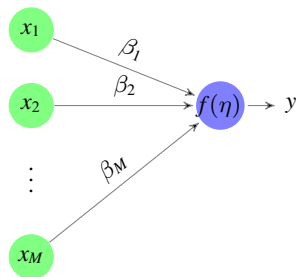
OVERVIEW OF DEEP NNs

Repetition of building blocks



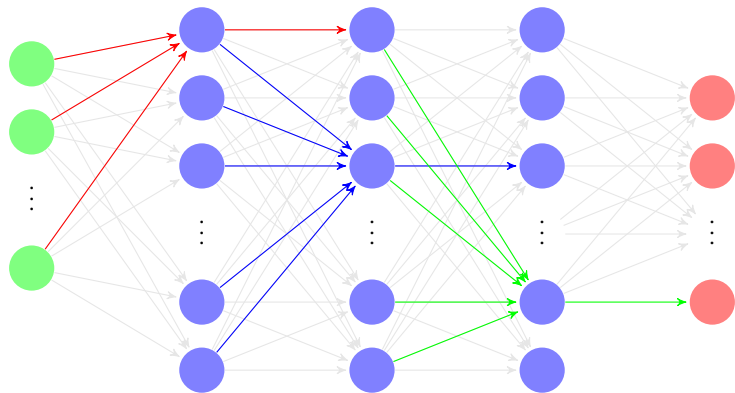
OVERVIEW OF DEEP NNs

Building block



- > Output: y
- > Linear predictor : $\eta = \sum_{m=1}^M \beta_m x_m = \beta^T \mathbf{x}$, coefficients β weights
- > Inverse Link function: $f(\eta)$ activation function

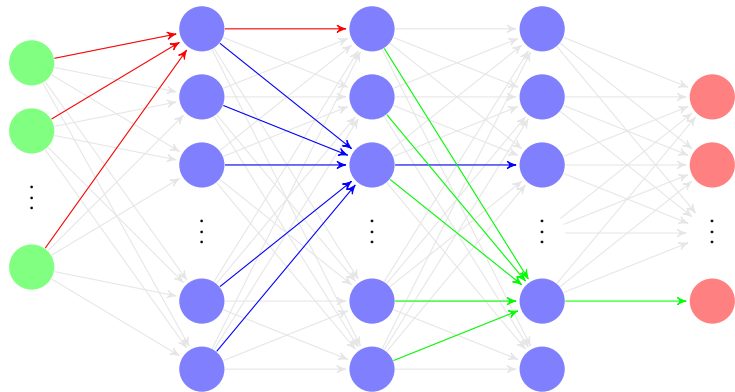
OVERVIEW OF DEEP NNs



In each layer l

the input is linearly composed $\rightarrow \eta_l = \beta_l^T \mathbf{x}_l$

OVERVIEW OF DEEP NNs

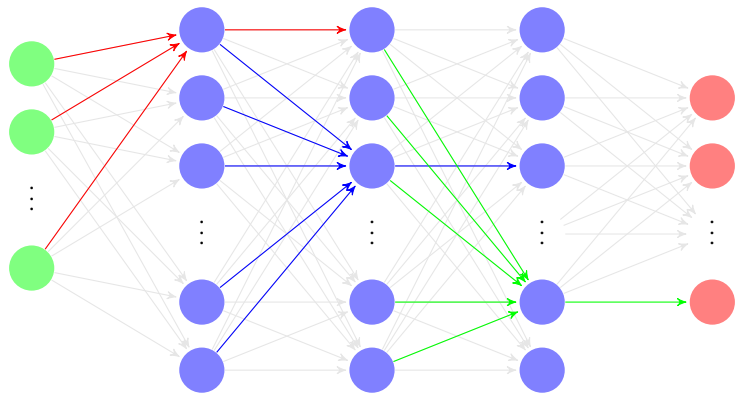


In each layer l

the input is linearly composed $\rightarrow \eta_l = \beta_l^T \mathbf{x}_l$

a non-linear function is applied $\rightarrow f_l(\eta_l)$

OVERVIEW OF DEEP NNs



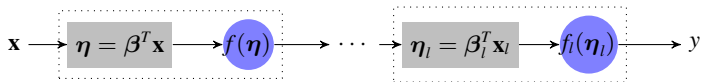
In each layer l

the input is linearly composed $\rightarrow \boldsymbol{\eta}_l = \boldsymbol{\beta}_l^T \mathbf{x}_l$

a non-linear function is applied $\rightarrow f_l(\boldsymbol{\eta}_l)$

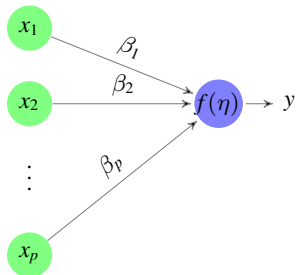
A Deep Neural Network \rightarrow construction of recursive building blocks

OVERVIEW OF DEEP NNs



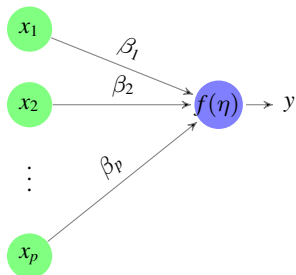
A Deep Neural Network can represent functions of increasing complexity.

DEEP NNs AS A PROBABILISTIC MODEL



$$\eta = \boldsymbol{\beta}^T \mathbf{x}$$

DEEP NNs AS A PROBABILISTIC MODEL



$$p(y|\mathbf{x}; \boldsymbol{\beta})$$

Parametric!

$$\eta = \boldsymbol{\beta}^T \mathbf{x}$$

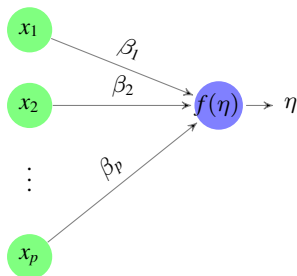
DEEP NNs AS A PROBABILISTIC MODEL



$$p(y|\mathbf{x}; \boldsymbol{\beta})$$

Parametric!

Example: Regression



$$\eta = \boldsymbol{\beta}^T \mathbf{x}$$

$$f(\eta) = \eta = \boldsymbol{\beta}^T \mathbf{x}$$
$$y = \boldsymbol{\beta}^T \mathbf{x} + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$p(y|\mathbf{x}, \boldsymbol{\beta}) = \mathcal{N}(\boldsymbol{\beta}^T \mathbf{x}, \sigma^2)$$

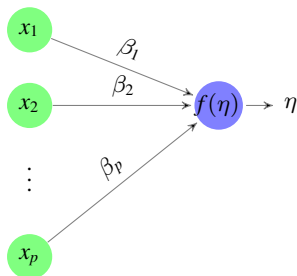
DEEP NNs AS A PROBABILISTIC MODEL



$$p(y|\mathbf{x}; \boldsymbol{\beta})$$

Parametric!

Example: Regression



$$\eta = \boldsymbol{\beta}^T \mathbf{x}$$

$$f(\eta) = \eta = \boldsymbol{\beta}^T \mathbf{x}$$
$$y = \boldsymbol{\beta}^T \mathbf{x} + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$p(y|\mathbf{x}, \boldsymbol{\beta}) = \mathcal{N}(\boldsymbol{\beta}^T \mathbf{x}, \sigma^2)$$

Train the network: learn $\boldsymbol{\beta}$?

Learn β

Maximum Likelihood Principal

Minimize the negative log-likelihood $\mathcal{J}(\beta) = -\log(p(\mathbf{Y}|\mathbf{x}, \beta))$

- $\beta_{MLE} : \arg \min_{\beta} J(\beta)$

DEEP NNS AS A PROBABILISTIC MODEL

Learn β

Maximum Likelihood Principal

Minimize the negative log-likelihood $\mathcal{J}(\beta) = -\log(p(\mathbf{Y}|\mathbf{x}, \beta))$

- $\beta_{MLE} : \arg \min_{\beta} J(\beta)$

Regression



$$p(y|\mathbf{x}, \beta) = \mathcal{N}(\beta^T \mathbf{x}, \sigma^2)$$

DEEP NNs AS A PROBABILISTIC MODEL

Learn β

Maximum Likelihood Principal

Minimize the negative log-likelihood $\mathcal{J}(\beta) = -\log(p(\mathbf{Y}|\mathbf{x}, \beta))$

- $\beta_{MLE} : \arg \min_{\beta} J(\beta)$

Regression



Squared Loss function

$$J(\beta) = \frac{1}{2} \sum_{i=1}^N \{y^{(i)} - \beta^T \mathbf{x}^{(i)}\}^2 + \text{const.}$$

$$p(y|\mathbf{x}, \beta) = \mathcal{N}(\beta^T \mathbf{x}, \sigma^2)$$

DEEP NNS AS A PROBABILISTIC MODEL

Learn β

Maximum Likelihood Principal

Minimize the negative log-likelihood $\mathcal{J}(\beta) = -\log(p(\mathbf{Y}|\mathbf{x}, \beta))$

- $\beta_{MLE} : \arg \min_{\beta} J(\beta)$

Regression



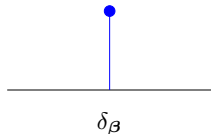
Squared Loss function

$$J(\beta) = \frac{1}{2} \sum_{i=1}^N \{y^{(i)} - \beta^T \mathbf{x}^{(i)}\}^2 + const.$$

$$p(y|\mathbf{x}, \beta) = \mathcal{N}(\beta^T \mathbf{x}, \sigma^2)$$

> Solve using Stochastic Gradient, back-propagation etc.

DEEP NNs AS A PROBABILISTIC MODEL

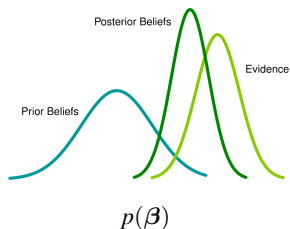
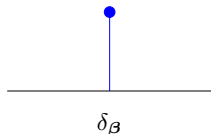


Maximum Likelihood Estimation

β is fixed but unknown

- Point estimates
- Overfitting

DEEP NNS AS A PROBABILISTIC MODEL



Maximum Likelihood Estimation

β is fixed but unknown

- Point estimates
- Overfitting

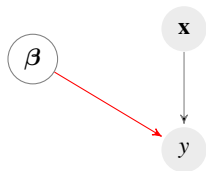
Bayesian Approach

β is unknown/uncertain

- + Account for uncertainty
 $p(\beta), p(\beta|\mathbf{y})$
- + Regularisation

BAYESIAN REASONING IN DEEP LEARNING

Maximum a Posteriori Estimation



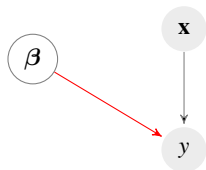
BAYESIAN REASONING IN DEEP LEARNING

Maximum a Posteriori Estimation

Distribution over the network weights β

$$\beta \sim p(\beta)$$

$$y|\mathbf{x}, \beta \sim p(y|\mathbf{x}, \beta)$$



BAYESIAN REASONING IN DEEP LEARNING

Maximum a Posteriori Estimation

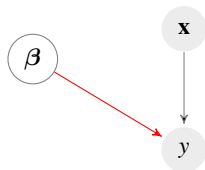
Distribution over the network weights β

$$\beta \sim p(\beta)$$

$$y|\mathbf{x}, \beta \sim p(y|\mathbf{x}, \beta)$$

Update the belief over $\beta \rightarrow$ Bayes' rule

$$p(\beta|\mathbf{Y}) \propto \overset{\text{Likelihood}}{p(\mathbf{Y}|\beta, \mathbf{X})} \overset{\text{prior}}{p(\beta)}$$



BAYESIAN REASONING IN DEEP LEARNING

Maximum a Posteriori Estimation

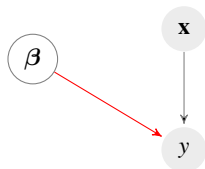
Distribution over the network weights β

$$\beta \sim p(\beta)$$

$$y|\mathbf{x}, \beta \sim p(y|\mathbf{x}, \beta)$$

Update the belief over $\beta \rightarrow$ Bayes' rule

$$p(\beta|\mathbf{Y}) \propto \overset{\text{Likelihood}}{p(\mathbf{Y}|\beta, \mathbf{X})} \overset{\text{prior}}{p(\beta)}$$



Maximum a Posteriori

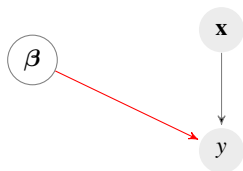
Find β_{MAP} that maximizes $\log p(\beta|\mathbf{Y})$

Find β_{MAP} that minimizes $\mathcal{J}(\beta) = -\log p(\beta|\mathbf{Y})$

- $\beta_{MAP} : \arg \min_{\beta} J(\beta)$

BAYESIAN REASONING IN DEEP LEARNING

Regression: Maximum a Posteriori



$$y = \beta^T \mathbf{x} + \epsilon$$

prior
over
the
weights

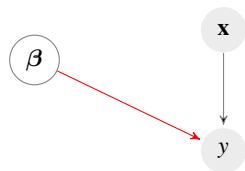
$$p(\beta) = \mathcal{N}(\mathbf{m}_\beta, \lambda^{-1} \mathbb{I})$$

multivariate
Gaussian

$$p(y|\mathbf{x}, \beta) = \mathcal{N}(\beta^T \mathbf{x}, \sigma^2)$$

BAYESIAN REASONING IN DEEP LEARNING

Regression: Maximum a Posteriori



Loss function: $\mathcal{J}(\beta) = -\log p(\beta|\mathbf{Y})$

$$y = \beta^T \mathbf{x} + \epsilon$$

prior
over
the
weights

$$p(\beta) = \mathcal{N}(\mathbf{m}_\beta, \lambda^{-1} \mathbb{I})$$

$$p(y|\mathbf{x}, \beta) = \mathcal{N}(\beta^T \mathbf{x}, \sigma^2)$$

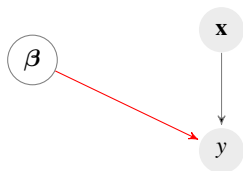
multivariate
Gaussian

$$\mathcal{J}(\beta) = \frac{1}{2} \sum_{i=1}^N \{y^{(i)} - \beta^T \mathbf{x}^{(i)}\}^2 + \frac{\lambda}{2} \beta^T \beta + \text{const.}$$

weight
decay -
regular-
ization

BAYESIAN REASONING IN DEEP LEARNING

Regression: Maximum a Posteriori



Loss function: $\mathcal{J}(\beta) = -\log p(\beta|\mathbf{Y})$

$$y = \beta^T \mathbf{x} + \epsilon$$

prior
over
the
weights

$$p(\beta) = \mathcal{N}(\mathbf{m}_\beta, \lambda^{-1} \mathbb{I})$$

$$p(y|\mathbf{x}, \beta) = \mathcal{N}(\beta^T \mathbf{x}, \sigma^2)$$

multivariate
Gaussian

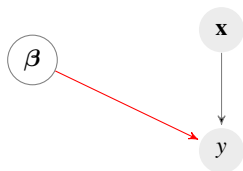
weight
decay -
regular-
ization

$$\mathcal{J}(\beta) = \frac{1}{2} \sum_{i=1}^N \{y^{(i)} - \beta^T \mathbf{x}^{(i)}\}^2 + \frac{\lambda}{2} \beta^T \beta + \text{const.}$$

The posterior distribution over the weights allows for a natural regularizer to arise! *Bayesian side effect*

BAYESIAN REASONING IN DEEP LEARNING

Regression: Maximum a Posteriori



Loss function: $\mathcal{J}(\beta) = -\log p(\beta|Y)$

$$y = \beta^T \mathbf{x} + \epsilon$$

prior
over
the
weights

$$p(\beta) = \mathcal{N}(\mathbf{m}_\beta, \lambda^{-1} \mathbb{I})$$

$$p(y|\mathbf{x}, \beta) = \mathcal{N}(\beta^T \mathbf{x}, \sigma^2)$$

multivariate
Gaussian

weight
decay -
regular-
ization

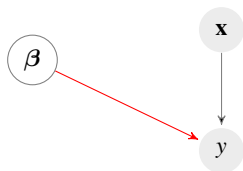
$$\mathcal{J}(\beta) = \frac{1}{2} \sum_{i=1}^N \{y^{(i)} - \beta^T \mathbf{x}^{(i)}\}^2 + \frac{\lambda}{2} \beta^T \beta + \text{const.}$$

The posterior distribution over the weights allows for a natural regularizer to arise! *Bayesian side effect*

$$\beta_{MAP} : \arg \min_{\beta} J(\beta)$$

BAYESIAN REASONING IN DEEP LEARNING

Regression: Maximum a Posteriori



Loss function: $\mathcal{J}(\beta) = -\log p(\beta|\mathbf{Y})$

$$y = \beta^T \mathbf{x} + \epsilon$$

prior
over
the
weights

multivariate
Gaussian

$$p(\beta) = \mathcal{N}(\mathbf{m}_\beta, \lambda^{-1} \mathbb{I})$$

$$p(y|\mathbf{x}, \beta) = \mathcal{N}(\beta^T \mathbf{x}, \sigma^2)$$

weight
decay -
regular-
ization

$$\mathcal{J}(\beta) = \frac{1}{2} \sum_{i=1}^N \{y^{(i)} - \beta^T \mathbf{x}^{(i)}\}^2 + \frac{\lambda}{2} \beta^T \beta + \text{const.}$$

The posterior distribution over the weights allows for a natural regularizer to arise! *Bayesian side effect*

$$\beta_{MAP} : \arg \min_{\beta} J(\beta)$$

Predictions:

$$y^* | \mathbf{Y}, \mathbf{x}^* \sim \mathcal{N}(\beta_{MAP}^T \mathbf{x}^*, \sigma^2)$$

What we have seen so far...

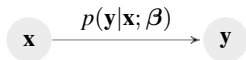
> Being Bayesian over the weights β , i.e. putting a prior distribution, allows for a better loss function in learning. **natural regularization**

BAYESIAN REASONING IN DEEP LEARNING

What we have seen so far...

> Being Bayesian over the weights β , i.e. putting a prior distribution, allows for a better loss function in learning. **natural regularization**

> (D)NNs **don't** handle missing data



$$\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, i = 1, \dots, N \text{ train data}$$
$$\mathbf{y}^* : p(\mathbf{y}^*|\mathbf{x}^*; \beta)$$

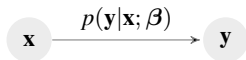
BAYESIAN REASONING IN DEEP LEARNING

What we have seen so far...

> Being Bayesian over the weights β , i.e. putting a prior distribution, allows for a better loss function in learning. **natural regularization**

> (D)NNs **don't** handle missing data

> Need for $p(\mathbf{x}); \mathbf{x} \sim p(\mathbf{x})$



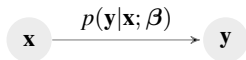
$$\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, i = 1, \dots, N \text{ train data}$$
$$\mathbf{y}^* : p(\mathbf{y}^*|\mathbf{x}^*; \beta)$$

$$\mathbf{x}^{(i)}?, \mathbf{y}^{(i)}$$

BAYESIAN REASONING IN DEEP LEARNING

What we have seen so far...

- > Being Bayesian over the weights β , i.e. putting a prior distribution, allows for a better loss function in learning. **natural regularization**
- > (D)NNs **don't** handle missing data
 - > Need for $p(\mathbf{x}); \mathbf{x} \sim p(\mathbf{x})$
- > (D)NNs **don't** generate new examples



$$\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, i = 1, \dots, N \text{ train data}$$
$$\mathbf{y}^* : p(\mathbf{y}^*|\mathbf{x}^*; \beta)$$

$$\mathbf{x}^{(i)} ?, \mathbf{y}^{(i)}$$

BAYESIAN REASONING IN DEEP LEARNING

What we have seen so far...

> Being Bayesian over the weights β , i.e. putting a prior distribution, allows for a better loss function in learning. **natural regularization**

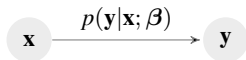
> (D)NNs **don't** handle missing data

> Need for $p(\mathbf{x}); \mathbf{x} \sim p(\mathbf{x})$

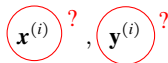
> (D)NNs **don't** generate new examples

> $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$

- draw $\mathbf{x}^{(i)} \sim p(\mathbf{x}^{(i)})$
- draw $\mathbf{y}^{(i)} \sim p(\mathbf{y}^{(i)}|\mathbf{x}^{(i)})$



$\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, i = 1, \dots, N$ train data
 $\mathbf{y}^* : p(\mathbf{y}^*|\mathbf{x}^*; \beta)$



BAYESIAN REASONING IN DEEP LEARNING

What we have seen so far...

> Being Bayesian over the weights β , i.e. putting a prior distribution, allows for a better loss function in learning. **natural regularization**

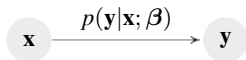
> (D)NNs **don't** handle missing data

> Need for $p(\mathbf{x}); \mathbf{x} \sim p(\mathbf{x})$

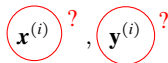
> (D)NNs **don't** generate new examples

> $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$

- draw $\mathbf{x}^{(i)} \sim p(\mathbf{x}^{(i)})$
- draw $\mathbf{y}^{(i)} \sim p(\mathbf{y}^{(i)}|\mathbf{x}^{(i)})$



$\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, i = 1, \dots, N$ train data
 $\mathbf{y}^* : p(\mathbf{y}^*|\mathbf{x}^*; \beta)$

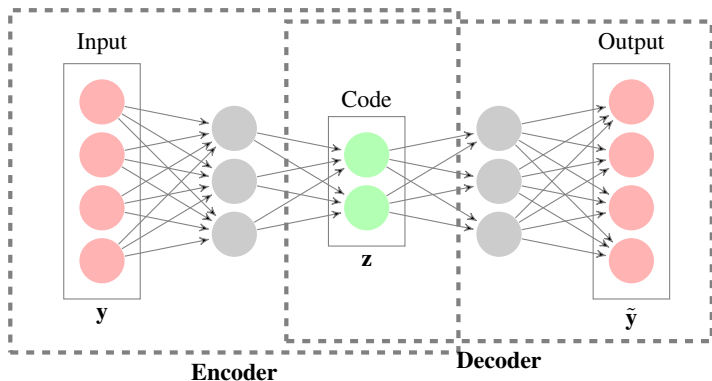


Need for a fully probabilistic model

Autoencoder as a fully probabilistic Neural Network

BAYESIAN REASONING IN DEEP LEARNING

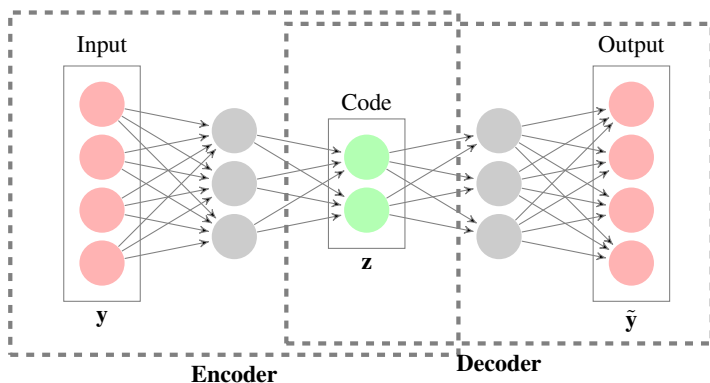
Autoencoder



- > Two parts:
 - > Encoder: "encodes" the input to hidden code (representation)
 - > Decoder: "decodes" the input back
- > Minimize reconstruction error $\mathcal{L} = \|y - \tilde{y}\|_2^2$

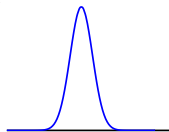
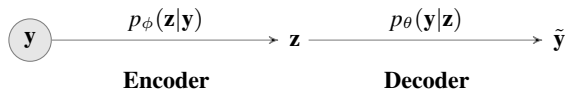
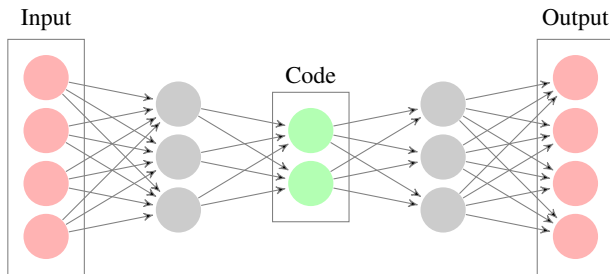
BAYESIAN REASONING IN DEEP LEARNING

Autoencoder



Unsupervised method: Experiences data \mathbf{y} and attempts to learn the distribution over \mathbf{y} , i.e. $p(\mathbf{y}|\mathbf{z})$

BAYESIAN REASONING IN DEEP LEARNING



$\{\theta, \phi\}$ NN's weights

BAYESIAN REASONING IN DEEP LEARNING

Decoder \rightarrow Generative Model



Latent variable
model

$$\mathbf{z} \sim p(\mathbf{z})$$

$$\tilde{\mathbf{y}} \sim p_{\theta}(\mathbf{y}|\mathbf{z})$$

BAYESIAN REASONING IN DEEP LEARNING

Encoder \rightarrow Inference
Mechanism



Learn $p_{\phi}(\mathbf{z}|\mathbf{y}) \rightarrow \phi$

Learn θ

Decoder \rightarrow Generative Model



Latent variable
model

$\mathbf{z} \sim p(\mathbf{z})$

$\tilde{\mathbf{y}} \sim p_{\theta}(\mathbf{y}|\mathbf{z})$

BAYESIAN REASONING IN DEEP LEARNING

Encoder \rightarrow Inference
Mechanism



Learn $p_{\phi}(\mathbf{z}|\mathbf{y}) \rightarrow \phi$

Learn θ

Decoder \rightarrow Generative Model



Latent variable
model

$\mathbf{z} \sim p(\mathbf{z})$

$\tilde{\mathbf{y}} \sim p_{\theta}(\mathbf{y}|\mathbf{z})$

Inference: Learn ϕ, θ weights!

Variational Inference

$$p_{\phi}(\mathbf{z}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{y})}?$$

Cannot sample from $p_{\phi}(\mathbf{z}|\mathbf{y}) \rightarrow$ cannot construct it using the NNs structure (weights ϕ)

Variational Inference

Main Principle

approximate distribution

$$q_{\phi}(\mathbf{z}) \approx p(\mathbf{z}|\mathbf{y})$$

$$p_{\phi}(\mathbf{z}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{y})}?$$

Cannot sample from $p_{\phi}(\mathbf{z}|\mathbf{y}) \rightarrow$ cannot construct it using the NNs structure (weights ϕ)

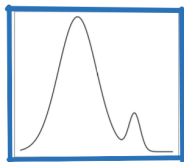
BAYESIAN REASONING IN DEEP LEARNING

Variational Inference

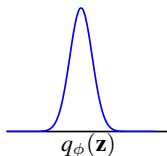
Main Principle

approximate distribution

$$q_{\phi}(\mathbf{z}) \approx p(\mathbf{z}|\mathbf{y})$$

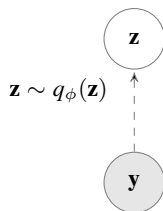


True posterior



$$p_{\phi}(\mathbf{z}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{y})}?$$

Cannot sample from $p_{\phi}(\mathbf{z}|\mathbf{y}) \rightarrow$ cannot construct it using the NNs structure (weights ϕ)



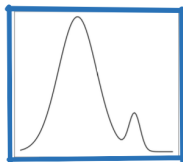
BAYESIAN REASONING IN DEEP LEARNING

Variational Inference

Main Principle

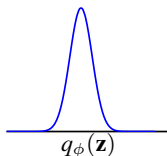
approximate distribution

$$q_{\phi}(\mathbf{z}) \approx p(\mathbf{z}|\mathbf{y})$$



True posterior

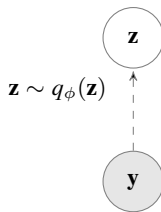
$KL[q_{\phi}(\mathbf{z})||p(\mathbf{z}|\mathbf{y})]$



$$p_{\phi}(\mathbf{z}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{y})}?$$

Cannot sample from $p_{\phi}(\mathbf{z}|\mathbf{y}) \rightarrow$ cannot construct it using the NNs structure (weights ϕ)

KL measures the information lost when using q to approximate p



Variational Inference

Cost function

A metric for how well $\{\phi, \theta\}$ capture the data generating $p_\theta(\mathbf{y}|\mathbf{z})$ & latent distributions $p_\phi(\mathbf{z}|\mathbf{y})$

Minimize

$$\mathcal{L} = -\log p(\mathbf{y}|\theta, \phi)$$

BAYESIAN REASONING IN DEEP LEARNING

Variational Inference

Cost function

A metric for how well $\{\phi, \theta\}$ capture the data generating $p_\theta(\mathbf{y}|\mathbf{z})$ & latent distributions $p_\phi(\mathbf{z}|\mathbf{y})$

Minimize

$$\mathcal{L} = -\log p(\mathbf{y}|\theta, \phi)$$

Reconstruction cost:

Expected log-likelihood measures how well samples from $q_\phi(\mathbf{z})$ are able to explain the data \mathbf{y} .

Penalty:

This divergence measures how much information is lost (in units of nats) when using $q_\phi(\mathbf{z})$ to represent $p(\mathbf{z}|\mathbf{y})$

Minimize \mathcal{F}

$$\mathcal{F}(\phi, \theta; \mathbf{y}) = \underbrace{-\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{y})}[\log p_\theta(\mathbf{y}|\mathbf{z})]}_{\text{Reconstruction}} + \underbrace{KL[q_\phi(\mathbf{z})||p(\mathbf{z}|\mathbf{y})]}_{\text{Penalty}}$$

BAYESIAN REASONING IN DEEP LEARNING

Variational Inference

Cost function

A metric for how well $\{\phi, \theta\}$ capture the data generating $p_\theta(\mathbf{y}|\mathbf{z})$ & latent distributions $p_\phi(\mathbf{z}|\mathbf{y})$

Reconstruction cost:

Expected log-likelihood measures how well samples from $q_\phi(\mathbf{z})$ are able to explain the data \mathbf{y} .

Penalty:

This divergence measures how much information is lost (in units of nats) when using $q_\phi(\mathbf{z})$ to represent $p(\mathbf{z}|\mathbf{y})$

Minimize \mathcal{F}

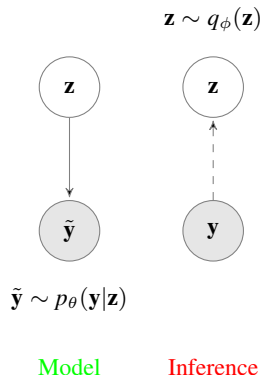
$$\mathcal{F}(\phi, \theta; \mathbf{y}) = \underbrace{-\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{y})}[\log p_\theta(\mathbf{y}|\mathbf{z})]}_{\text{Reconstruction}} + \underbrace{KL[q_\phi(\mathbf{z})||p(\mathbf{z}|\mathbf{y})]}_{\text{Penalty}}$$

Penalty is derived from your model and does not need to be designed

BAYESIAN REASONING IN DEEP LEARNING

Bayesian & DL marriage - What have we gained?

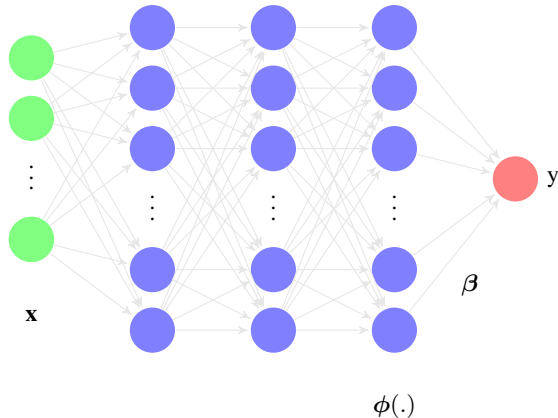
- ++ Principled inference approach - Bayesian
 - Penalty as intrinsic part of the model!
- ++ Encode uncertainty
- ++ Impute missing data



Deep Learning as a kernel method

- > **So far...** Predictions made using point estimates β (weights) **parametric**
- > **Now...** Store the entire training set in order to make predictions \rightarrow *kernel approach* **non-parametric**

DEEP LEARNING AS A KERNEL METHOD



Regression example

$$y = \beta^T \phi(\mathbf{x}) = \sum_{m=1}^M \beta_m \phi_m(\mathbf{x})$$

Parametric approach

$$J(\boldsymbol{\beta}) = \frac{1}{2} \sum_{n=1}^N \{\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n) - y_n\}^2 + \frac{\lambda}{2} \boldsymbol{\beta}^T \boldsymbol{\beta}$$

- > MLE training \rightarrow point estimates of $\boldsymbol{\beta}$

$$\nabla J(\boldsymbol{\beta}) = 0 \rightarrow \boldsymbol{\beta}_{\text{MLE}} = -\frac{1}{\lambda} \sum_{n=1}^N \{\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n) - y_n\} \boldsymbol{\phi}(\mathbf{x}_n)$$

- > Predictions:

$$y^* = \boldsymbol{\beta}_{\text{MLE}}^T \boldsymbol{\phi}(\mathbf{x}^*)$$

Predictions purely based on the learned parameter- all the information from the data is transferred into a set of parameters

Parametric approach

$$J(\boldsymbol{\beta}) = \frac{1}{2} \sum_{n=1}^N \{\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n) - y_n\}^2 + \frac{\lambda}{2} \boldsymbol{\beta}^T \boldsymbol{\beta}$$

> MLE training \rightarrow point estimates of $\boldsymbol{\beta}$

$$\nabla J(\boldsymbol{\beta}) = 0 \rightarrow \boldsymbol{\beta}_{\text{MLE}} = -\frac{1}{\lambda} \sum_{n=1}^N \{\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n) - y_n\} \boldsymbol{\phi}(\mathbf{x}_n)$$

> Predictions:

$$y^* = \boldsymbol{\beta}_{\text{MLE}}^T \boldsymbol{\phi}(\mathbf{x}^*)$$

Predictions purely based on the learned parameter- all the information from the data is transferred into a set of parameters

parametric!

DEEP LEARNING AS A KERNEL METHOD

Change of scenery...let the data speak

$$J(\boldsymbol{\beta}) = \frac{1}{2} \sum_{n=1}^N \{\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n) - y_n\}^2 + \frac{\lambda}{2} \boldsymbol{\beta}^T \boldsymbol{\beta}$$

$$\nabla J(\boldsymbol{\beta}) = 0 \rightarrow$$

$$\boldsymbol{\beta}_{\text{MLE}} = -\frac{1}{\lambda} \sum_{n=1}^N \{\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n) - y_n\} \boldsymbol{\phi}(\mathbf{x}_n)$$

$$= \sum_{n=1}^N \alpha_n \boldsymbol{\phi}(\mathbf{x}_n) = \begin{matrix} \text{design} \\ \text{matrix} \\ M \times N \end{matrix} \boldsymbol{\Phi}^T \quad \begin{matrix} N \times 1 \\ \boldsymbol{\alpha} \end{matrix}$$

$$\alpha_n = -\frac{1}{\lambda} (\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n) - y_n)$$

DEEP LEARNING AS A KERNEL METHOD

Change of scenery...let the data speak

$$J(\boldsymbol{\beta}) = \frac{1}{2} \sum_{n=1}^N \{\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n) - y_n\}^2 + \frac{\lambda}{2} \boldsymbol{\beta}^T \boldsymbol{\beta}$$

dual representation

$$\nabla J(\boldsymbol{\beta}) = 0 \rightarrow$$

$$\boldsymbol{\beta}_{\text{MLE}} = -\frac{1}{\lambda} \sum_{n=1}^N \{\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n) - y_n\} \boldsymbol{\phi}(\mathbf{x}_n)$$

$$= \sum_{n=1}^N \alpha_n \boldsymbol{\phi}(\mathbf{x}_n) = \begin{matrix} \text{design} \\ \text{matrix} \\ M \times N \end{matrix} \boldsymbol{\Phi}^T \begin{matrix} N \times 1 \\ \boldsymbol{\alpha} \end{matrix}$$

$$\alpha_n = -\frac{1}{\lambda} (\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n) - y_n)$$

$$\nabla J(\boldsymbol{\alpha}) = 0 \rightarrow$$

Gram
matrix
 $N \times N$

$$\boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{y}$$

$$\kappa_{lm} = \langle \boldsymbol{\phi}(\mathbf{x}_l), \boldsymbol{\phi}(\mathbf{x}_m) \rangle = \boxed{k(\mathbf{x}_l, \mathbf{x}_m)}$$

kernel function

Why?

$$\beta_{\text{MLE}} = -\frac{1}{\lambda} \sum_{n=1}^N \{\beta^T \phi(\mathbf{x}_n) - y_n\} \phi(\mathbf{x}_n)$$

Why?

$$\boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{y}$$

$$\boldsymbol{\beta}_{\text{MLE}} = -\frac{1}{\lambda} \sum_{n=1}^N \{\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n) - y_n\} \boldsymbol{\phi}(\mathbf{x}_n)$$

DEEP LEARNING AS A KERNEL METHOD

Why?

$$\beta_{\text{MLE}} = -\frac{1}{\lambda} \sum_{n=1}^N \{\beta^T \phi(\mathbf{x}_n) - y_n\} \phi(\mathbf{x}_n)$$

$$\alpha = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{y}$$

Predictions:

point
estimate-
no
mem-
ory of
data

$$y^* = \beta_{\text{MLE}}^T \phi(\mathbf{x}^*)$$

DEEP LEARNING AS A KERNEL METHOD

Why?

$$\beta_{\text{MLE}} = -\frac{1}{\lambda} \sum_{n=1}^N \{\beta^T \phi(\mathbf{x}_n) - y_n\} \phi(\mathbf{x}_n)$$

$$\alpha = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{y}$$

Predictions:

point
estimate-
no
mem-
ory of
data

$$y^* = \beta_{\text{MLE}}^T \phi(\mathbf{x}^*)$$

all data

$$y^* = k(\mathbf{X}, \mathbf{x}^*)^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

DEEP LEARNING AS A KERNEL METHOD

Why?

$$\beta_{\text{MLE}} = -\frac{1}{\lambda} \sum_{n=1}^N \{\beta^T \phi(\mathbf{x}_n) - y_n\} \phi(\mathbf{x}_n)$$

$$\alpha = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{y}$$

Predictions:

point estimate-
no memory of data

$$y^* = \beta_{\text{MLE}}^T \phi(\mathbf{x}^*)$$

all data

$$y^* = k(\mathbf{X}, \mathbf{x}^*)^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

parametric \rightarrow non parametric

DEEP LEARNING AS A KERNEL METHOD

Why?

$$\beta_{\text{MLE}} = -\frac{1}{\lambda} \sum_{n=1}^N \{\beta^T \phi(\mathbf{x}_n) - y_n\} \phi(\mathbf{x}_n)$$

$$\alpha = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{y}$$

Predictions:

point estimate-
no memory of data

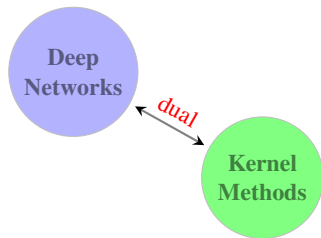
$$y^* = \beta_{\text{MLE}}^T \phi(\mathbf{x}^*)$$

all data

$$y^* = k(\mathbf{X}, \mathbf{x}^*)^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

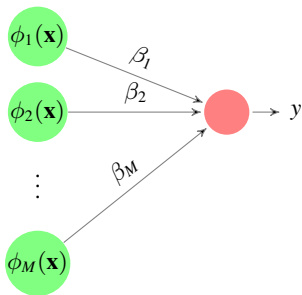
parametric \rightarrow non parametric
Memory in β_{MLE} \rightarrow Memory by actual storing all the data

DEEP LEARNING AS A KERNEL METHOD



DEEP LEARNING AS A KERNEL METHOD

Step back...

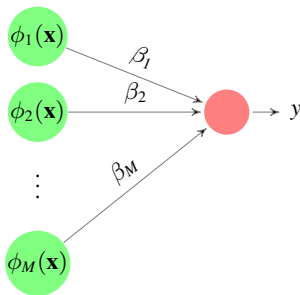


$$y = \boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}) = \sum_m^M \beta_m \phi_m(\mathbf{x})$$
$$y = f(\mathbf{x})$$

TASK: learn the best regression function $f()$

DEEP LEARNING AS A KERNEL METHOD

Step back...



$$y = \boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}) = \sum_m^M \beta_m \phi_m(\mathbf{x})$$
$$y = f(\mathbf{x})$$

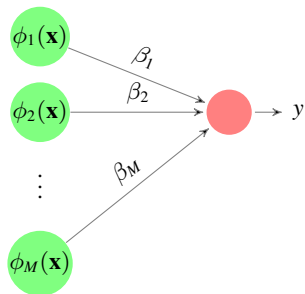
TASK: learn the best regression function $f()$

What about a bit more of Bayesian reasoning?

$$\mathbf{f} = [f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n)] \sim p(\mathbf{f})$$

Consider \mathbf{f} a multivariate variable- Apply distribution!

DEEP LEARNING AS A KERNEL METHOD



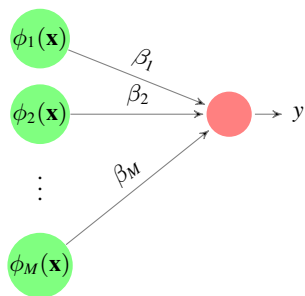
- Number of hidden units $\phi_m(\mathbf{x})$ to **infinity**, $M \rightarrow \infty$
- $\beta \sim \mathcal{N}(\mathbf{m}, \Sigma)$

$$y = \beta^T \phi(\mathbf{x}) = \sum_m^M \beta_m \phi_m(\mathbf{x})$$
$$y = f(\mathbf{x})$$

$$\mathbf{f} = [f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n)]$$

$$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}), \quad K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$$

DEEP LEARNING AS A KERNEL METHOD



- Number of hidden units $\phi_m(\mathbf{x})$ to **infinity**, $M \rightarrow \infty$
- $\beta \sim \mathcal{N}(\mathbf{m}, \Sigma)$

$$y = \beta^T \phi(\mathbf{x}) = \sum_m^M \beta_m \phi_m(\mathbf{x})$$
$$y = f(\mathbf{x})$$

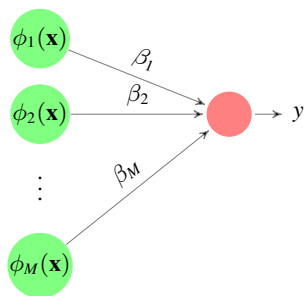
$$\mathbf{f} = [f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n)]$$

$$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}), \quad K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$$

When $N \rightarrow \infty$, sample from a **Gaussian process**

$$\mathbf{f} \sim \text{GP}(\mathbf{0}, \mathbf{K})$$

DEEP LEARNING AS A KERNEL METHOD



- Number of hidden units $\phi_m(\mathbf{x})$ to **infinity**, $M \rightarrow \infty$
- $\beta \sim \mathcal{N}(\mathbf{m}, \Sigma)$

$$y = \beta^T \phi(\mathbf{x}) = \sum_m^M \beta_m \phi_m(\mathbf{x})$$
$$y = f(\mathbf{x})$$

$$\mathbf{f} = [f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n)]$$

$$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}), \quad K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$$

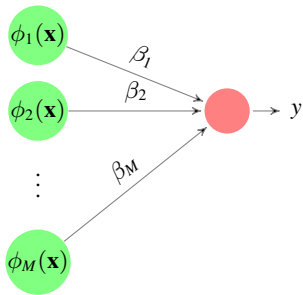
When $N \rightarrow \infty$, sample from a **Gaussian process**

$$\mathbf{f} \sim \text{GP}(\mathbf{0}, \mathbf{K})$$

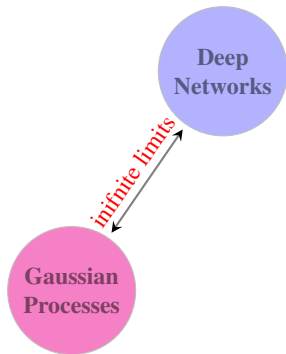
Predict $y^* = f(\mathbf{x}^*)$

$$p(f^* | \mathbf{X}, \mathbf{y}, \mathbf{x}^*) = \int p(y^* | \mathbf{x}^*, \mathbf{f}, \mathbf{X}, \mathbf{y}) p(\mathbf{f} | \mathbf{X}, \mathbf{y}) d\mathbf{f}$$

DEEP LEARNING AS A KERNEL METHOD



Number of hidden units $\phi_m(\mathbf{x})$ to
infinity, $M \rightarrow \infty$



Thank you!