Probabilistic Programming with Deep Neural Networks

Andrés Masegosa

Department of Mathematics University of Almería Spain

https://andresmasegosa.github.io/ https://pgmlab.github.io/

This material jointly made with Thomas D. Nielsen (AAU).

Day 1: Introduction to probabilistic programming languages

- Why do we need PPLs?
- Probabilistic programming in Pyro
- Hand-on exercises:
 - Probability Distributions in Pyro.
 - Probabilistic Models in Pyro.

Day 2: Probabilistic Models with Deep Neural Networks

- Uncertainty in Machine Learning
- Variational Inference
- Supervised/Unsupervised Learning
- Hand-on exercises
 - Bayesian Neural Networks
 - Variational Auto-encoders

What is a PPL?

What is a probabilistic model?





VQ-VAE-2



Probabilistic Programming Language (PPL)

• An attempt to unify probabilistic modeling and general programming languages.



Probabilistic Programming Language (PPL)

- An attempt to unify probabilistic modeling and general programming languages.
- A programming paradigm to define general probabilistic models (mixing deterministic + stochastic functions).



Probabilistic Programming Language (PPL)

- An attempt to unify probabilistic modeling and general programming languages.
- A programming paradigm to define general probabilistic models (mixing deterministic + stochastic functions).
- Make probabilistic modeling more applicable and powerful.

Why PPLs?



Why do we need PPLs?



Why do we need PPLs?

• Reason 1: Try to democratize the development of AI systems.



Why do we need PPLs?

- Reason 1: Try to democratize the development of AI systems.
- Reason 2: Try to make AI systems safer.

Reason 1: Try to democratize the development of AI systems.



DARPA's Fund Call for PPLs in Artificial Intelligence.



DARPA's Fund Call for PPLs in Artificial Intelligence.

The development of machine learning systems requires enormous efforts.

Data Science



Copyright © 2014 by Steven Geringer Raleigh, NC. Permission is granted to use, distribute, or modify this image, provided that this copyright notice remains intact.

The development of machine learning systems requires enormous efforts.

• It requires of highly qualified experts.

Why PPLs?



The development of machine learning applications requires enormous effort.

- It is necessary to have highly qualified experts.
- It is difficult to find the ML model most suitable for an application.

Hidden Technical Debt in Machine Learning Systems



The development of machine learning applications requires enormous effort.

- It is necessary to have highly qualified experts.
- It is difficult to find the ML model most suitable for an application.
- Programming a ML model is a complex task where many problems are intermingled.

Wanted: Artificial intelligence experts

In artificial intelligence, job openings are rising faster than job seekers.



Consequences:

• Shortage of AI experts (and high salaries).

Wanted: Artificial intelligence experts

In artificial intelligence, job openings are rising faster than job seekers.



Consequences:

- Shortage of AI experts (and high salaries).
- Only big corporations have the resources for developing ML systems.





• People used to program in low-level programming languages.



- People used to program in low-level programming languages.
- Programming was complex and demand high-expertise.



- People used to program in low-level programming languages.
- Programming was complex and demand high-expertise.
- Focus on application and low-level hardware details.





• Programmers focused on the applications.



- Programmers focused on the applications.
- Hardware Experts focused on compilers.



- Programmers focused on the applications.
- Hardware Experts focused on compilers.
- High gains in productivity.



- Programmers focused on the applications.
- Hardware Experts focused on compilers.
- High gains in productivity.
- "Democratization" of the software development.



• Stacked architecture



• Stacked architecture

• Different Domain Experts will code their models using the same language.



- Stacked architecture
- Different Domain Experts will code their models using the same language.
- ML experts will focus on the development of new ML solvers.



- Stacked architecture
- Different Domain Experts will code their models using the same language.
- ML experts will focus on the development of new ML solvers.
- Compile experts will focus on running these ML solvers on specialized hardware.



Benefits of PPLs:

• Simplify machine learning model code.



Benefits of PPLs:

- Simplify machine learning model code.
- Reduce development time and cost to encourage experimentation.



Benefits of PPLs:

- Simplify machine learning model code.
- Reduce development time and cost to encourage experimentation.
- Facilitate the construction of more sophisticated models.


Benefits of PPLs:

- Simplify machine learning model code.
- Reduce development time and cost to encourage experimentation.
- Facilitate the construction of more sophisticated models.
- Reduce the necessary level of expertise.



Benefits of PPLs:

- Simplify machine learning model code.
- Reduce development time and cost to encourage experimentation.
- Facilitate the construction of more sophisticated models.
- Reduce the necessary level of expertise.
- "Democratization" of the development of ML systems.

Reason 2: Try to make AI systems safer



Deep Learning Based AI Systems :

- Issue 1: Hard to interpret.
- Issue 2: No possible to know how sure they are in a particular prediction.
- Enormously limit the application of AI to many real life problems.



AI in Health Care:

- Patients need to know why they are prescribed some treatment.
- An doctor can supervise the machine (humans in the loop).
- Extensible to any safe-critical system.



AI in Automated Systems:

- The system should detect when it is in a completely new situation.
- Let a human take the control.











DARPA's Fund Call for XAI projects.



PPL based Systems :

- Addressing Issue 1: PPLs provide transparent model description.
- Addressing Issue 2: PPLs provide **uncertainty estimation** both in models and predictions.

Brief Historical Review of PPLs



• Bugs, WinBugs, Jags, Figaro, etc.



- Bugs, WinBugs, Jags, Figaro, etc.
- **Turing-complete** probabilistic programming languages. (i.e. they can represent any computable probability distribution).



- Bugs, WinBugs, Jags, Figaro, etc.
- **Turing-complete** probabilistic programming languages. (i.e. they can represent any computable probability distribution).
- Inference engine based on Monte Carlo methods.



- Bugs, WinBugs, Jags, Figaro, etc.
- **Turing-complete** probabilistic programming languages. (i.e. they can represent any computable probability distribution).
- Inference engine based on Monte Carlo methods.
- They did not scale to large data samples/high-dimensional models.



• Infer.net, Factorie, Amidst, etc.



- Infer.net, Factorie, Amidst, etc.
- Inference engine based on **message passage algorithms** and/or variational inference methods.



- Infer.net, Factorie, Amidst, etc.
- Inference engine based on **message passage algorithms** and/or variational inference methods.
- They did scale to large data samples/high-dimensional models.



- Infer.net, Factorie, Amidst, etc.
- Inference engine based on **message passage algorithms** and/or variational inference methods.
- They did scale to large data samples/high-dimensional models.
- **Restricted** probabilistic model family (i.e. factor graphs, conjuage exponential family, etc.)



• TensorFlow Probability, Pyro, PyMC3, InferPy, etc.



- TensorFlow Probability, Pyro, PyMC3, InferPy, etc.
- Black Box Variational Inference and Hamiltonian Monte-Carlo.



- TensorFlow Probability, Pyro, PyMC3, InferPy, etc.
- Black Box Variational Inference and Hamiltonian Monte-Carlo.
- They did scale to large data samples/high-dimensional models.



- TensorFlow Probability, Pyro, PyMC3, InferPy, etc.
- Black Box Variational Inference and Hamiltonian Monte-Carlo.
- They did scale to large data samples/high-dimensional models.
- Turing-complete probabilistic programming languages.



- TensorFlow Probability, Pyro, PyMC3, InferPy, etc.
- Black Box Variational Inference and Hamiltonian Monte-Carlo.
- They did scale to large data samples/high-dimensional models.
- Turing-complete probabilistic programming languages.
- Allow the inclusion of deep neural networks.



- TensorFlow Probability, Pyro, PyMC3, InferPy, etc.
- Black Box Variational Inference and Hamiltonian Monte-Carlo.
- They did scale to large data samples/high-dimensional models.
- Turing-complete probabilistic programming languages.
- Allow the inclusion of deep neural networks.
- Rely on deep learning frameworks (TensorFlow, Pytorch, Theano, etc).

Pyro





• Developed by **UBER** (the car riding company).



- Developed by **UBER** (the car riding company).
- Focus on probabilistic models with deep neural networks.



- Developed by **UBER** (the car riding company).
- Focus on probabilistic models with deep neural networks.
- Rely on Pytorch (Deep Learning Framework).



- Developed by **UBER** (the car riding company).
- Focus on probabilistic models with deep neural networks.
- Rely on **Pytorch** (Deep Learning Framework).
- Enable GPU accelaration and distributed learning.



Demand Prediction with Pyro:

• Demand prediction is critical for user experience, resource allocation, etc.



Demand Prediction with Pyro:

- Demand prediction is critical for user experience, resource allocation, etc.
- LSTM powerful for time series modelling.



Demand Prediction with Pyro:

- Demand prediction is critical for user experience, resource allocation, etc.
- LSTM powerful for time series modelling.
- Prediction at special events is challenging: weather, population growth, etc.



Demand Prediction with Pyro:

- Demand prediction is critical for user experience, resource allocation, etc.
- LSTM powerful for time series modelling.
- Prediction at special events is challenging: weather, population growth, etc.
- Bayesian LSTM provides uncertainty estimation.

Pyro's Distributions


Pyro's distributions (http://docs.pyro.ai/en/stable/distributions.html) :

• Wide range of distributions: Normal, Beta, Cauchy, Dirichlet, Gumbel, Poisson, Pareto, etc.

M In [15]: sample = normal.sample()
sample
Out[15]: tensor(0.4908)

- Wide range of distributions: Normal, Beta, Cauchy, Dirichlet, Gumbel, Poisson, Pareto, etc.
- Samples from the distributions are Pytorch's Tensor objects (i.e. multidimensional arrays).

In [17]: sample = normal.sample(sample_shape=[3,4,5])
sample.shape

Out[17]: torch.Size([3, 4, 5])

- Wide range of distributions: Normal, Beta, Cauchy, Dirichlet, Gumbel, Poisson, Pareto, etc.
- Samples from the distributions are Pytorch's Tensor objects (i.e. multidimensional arrays).

In [19]: torch.sum(normal.log_prob(sample))

Out[19]: tensor(-85.1003)

- Wide range of distributions: Normal, Beta, Cauchy, Dirichlet, Gumbel, Poisson, Pareto, etc.
- Samples from the distributions are Pytorch's Tensor objects (i.e. multidimensional arrays).
- Operations, like log-likelihood, are defined over tensors (with GPU acceleration powered by Pytorch).

In [9]: normal = dist.Normal(torch.tensor([1.,2.,3.]),1.)
normal

Out[9]: Normal(loc: torch.Size([3]), scale: torch.Size([3]))

- Wide range of distributions: Normal, Beta, Cauchy, Dirichlet, Gumbel, Poisson, Pareto, etc.
- Samples from the distributions are Pytorch's Tensor objects (i.e. multidimensional arrays).
- Operations, like log-likelihood, are defined over tensors (with GPU acceleration powered by Pytorch).
- Multiple distributions can be embedded in single object (to define efficient vectorized operations).

M In [10]: normal.sample() Out[10]: tensor([2.0592, 2.4035, 3.1918])

- Wide range of distributions: Normal, Beta, Cauchy, Dirichlet, Gumbel, Poisson, Pareto, etc.
- Samples from the distributions are Pytorch's Tensor objects (i.e. multidimensional arrays).
- Operations, like log-likelihood, are defined over tensors (with GPU acceleration powered by Pytorch).
- Multiple distributions can be embedded in single object (to define efficient vectorized operations).

In [11]: normal.log_prob(normal.sample())

Out[11]: tensor([-0.9402, -1.2113, -2.3214])

- Wide range of distributions: Normal, Beta, Cauchy, Dirichlet, Gumbel, Poisson, Pareto, etc.
- Samples from the distributions are Pytorch's Tensor objects (i.e. multidimensional arrays).
- Operations, like log-likelihood, are defined over tensors (with GPU acceleration powered by Pytorch).
- Multiple distributions can be embedded in single object (to define efficient vectorized operations).

Open the Notebook and Play around

- Test that you have installed the basic packages.
- Test that you can run the first lines of code.
- Play a bit with the code in Section 1 of the notebook.

Day1/students_PPLs_Intro.ipynb

https://github.com/PGM-Lab/ASML-Tbilisi

Pyro's Models

```
M In [12]: def model():
    temp = pyro.sample('temp', dist.Normal(15.0, 2.0))
    return temp
print(model())
tensor(12.3926)
tensor(22.5272)
```

• A probabilistic model is defined as a stochastic function.

```
M In [12]: def model():
    temp = pyro.sample('temp', dist.Normal(15.0, 2.0))
    return temp
print(model())
tensor(12.3926)
tensor(22.5272)
```

- A probabilistic model is defined as a stochastic function.
- Each random variable is associated to a **primitive stochastic function** using the construct **pyro.sample(...)**.

```
M In [21]: def model():
    temp = pyro.sample('temp', dist.Normal(15.0, 2.0))
    sensor = pyro.sample('sensor', dist.Normal(temp, 1.0))
    return (temp, sensor)
    outl = model()
    outl
Out[21]: (tensor(15.9576), tensor(16.9907))
```

• A stochastic function can be defined as a composition of primitive stochastic functions.

```
M In [21]: def model():
    temp = pyro.sample('temp', dist.Normal(15.0, 2.0))
    sensor = pyro.sample('sensor', dist.Normal(temp, 1.0))
    return (temp, sensor)
    outl = model()
    outl
Out[21]: (tensor(15.9576), tensor(16.9907))
```

- A stochastic function can be defined as a composition of primitive stochastic functions.
- We define the joint probability distribution:

```
p(\texttt{sensor}, \texttt{temp}) = p(\texttt{sensor}|\texttt{temp})p(\texttt{temp})
```

Pyro's Inference



• We can introduce observations (e.g. sensor = 18.0).

```
M In [27]: #Run inference
svi(model,guide,obs, plot=True)
#Print results
print("P(Temperature|Sensor=18.0) = ")
print(dist.Normal(pyro.param("mean").item(), pyro.param("scale").item()))
print("")
```

- We can introduce observations (e.g. sensor = 18.0).
- We can query the posterior probability distribution:

$$p(\texttt{temp}|\texttt{sensor} = 18) = \frac{p(\texttt{sensor} = 18 | \texttt{temp}) p(\texttt{temp})}{\int p(\texttt{sensor} = 18 | \texttt{temp}) p(\texttt{temp}) d\texttt{temp}}$$

```
M In [27]: #Run inference
svi(model,guide,obs, plot=True)
#Print results
print("P(Temperature|Sensor=18.0) = ")
print(dist.Normal(pyro.param("mean").item(), pyro.param("scale").item()))
print("")
```

- We can introduce observations (e.g. sensor = 18.0).
- We can query the posterior probability distribution:

$$p(\texttt{temp}|\texttt{sensor} = 18) = \frac{p(\texttt{sensor} = 18 | \texttt{temp}) p(\texttt{temp})}{\int p(\texttt{sensor} = 18 | \texttt{temp}) p(\texttt{temp}) d\texttt{temp}}$$

• Guide is an auxiliary method needed for inference (more details in the coming sessions).



P(Temperature|Sensor=18.0) = Normal(loc: 17.39859390258789, scale: 0.9089401960372925)



P(Temperature|Sensor=18.0) = Normal(loc: 17.39859390258789, scale: 0.9089401960372925)

Details on the inference method will be given on the following sessions.



https://github.com/PGM-Lab/ASML-Tbilisi



• What if we have a **bunch of observations**, $\mathbf{s} = \{s_1, \dots, s_n\}$



- What if we have a **bunch of observations**, $\mathbf{s} = \{s_1, \dots, s_n\}$
- A random variable is created for each observation (using a for-loop).



• What if we do not know the average temperature?



- What if we do not know the average temperature?
- We can introduce a parameter using pyro.param construct.



• And learn the parameter with the same general inference algorithm.

$$\mu_t = \arg\max_{\mu} \ln p(s_1, \dots, s_n | \mu)$$

• Details about the inference algorithm will be given in the next sessions.



• What if we want to capture uncertainty about the estimation of the average temperature?



- What if we want to capture uncertainty about the estimation of the average temperature?
- We can model this parameter with a random variable.

```
In [162]: import time
    #Run inference
    start = time.time()
    svi(model, guide, obs, num_steps=1000)
    #Print results
    print("P(mean_temp|Sensor=[18., 18.7, 19.2, 17.8, 20.3, 22.4, 20.3, 21.2, 19.5, 20.1]) =")
    print(dist.Normal(pyro.param("mean").item(), pyro.param("scale").item()))
    print(t"")
    end = time.time()
    print(f"(end - start)) seconds")
    P(mean_temp|Sensor=[18., 18.7, 19.2, 17.8, 20.3, 22.4, 20.3, 21.2, 19.5, 20.1]) =
    Normal(loc: 19.199871063232422, scale: 0.6046891212463379)
    10.298431873321533 seconds
```

• And learn the distribution with the same general inference algorithm.

 $p(\mu_t|s_1,\ldots,s_{10})$

• Details about the inference algorithm will be given in the next sessions.



• And learn the distribution with the same general inference algorithm.

$$p(\mu_t|s_1,\ldots,s_{10})$$

• Details about the inference algorithm will be given in the next sessions.

f

Defining Conditional Independences in Pyro



Pyro's cond. independeces (http://pyro.ai/examples/svi_part_ii.html) :

• Sensor variables are **independent** given temperature mean, μ_t .

$$p(s_1, t_1, ..., s_{10}, t_{10} | \mu_t) = \prod_{i=1}^{10} p(s_i, t_i | \mu_t)$$



Pyro's cond. independeces (http://pyro.ai/examples/svi_part_ii.html) :

• Sensor variables are **independent** given temperature mean, μ_t .

$$p(s_1, t_1, ..., s_{10}, t_{10} | \mu_t) = \prod_{i=1}^{10} p(s_i, t_i | \mu_t)$$

• We can use Pyro's plate construct to introduce this independence.



```
2.81210994720459 seconds
```

Pyro's cond. independeces (http://pyro.ai/examples/svi_part_ii.html) :

- We get large gains in efficiency due to vectorized operations.
- Execution time without *plate* is over 10s.

- Exercise 2: The role of the number of observations in learning.
- Exercise 3: The role of the prior distribution in learning.

Day1/students_PPLs_Intro.ipynb

Toy Example: Ice-cream shop



Defining Machine Learning models with PPLs:

- We have an ice-cream shop and we record the ice-cream sales and the average temperature of the day.
- We know temperature affects the sales of ice-creams.
- We want to precisely find out how temperature affects ice-cream sales.
Ice-cream Shop Model:

• We have observations from temperature and sales.



Ice-cream Shop Model:

- We have observations from temperature and sales.
- Sales are modeled with a Poisson distribution.



Ice-cream Shop Model:

- We have observations from temperature and sales.
- Sales are modeled with a Poisson distribution.
- The rate of the Poisson linearly depends of the real temperature.

Pyro

```
M In [48]: #Run inference
           svi(model, guide, obs, num steps=1000)
           #Print results
           print("Posterior Temperature Mean")
           print(dist.Normal(pyro.param("mean").item(), pyro.param("scale").item()))
           print("")
           print("Posterior Alpha")
           print(dist.Normal(pyro.param("alpha mean").item(), pyro.param("alpha scale").item()))
           print("")
           print("Posterior Beta")
           print(dist.Normal(pyro.param("beta mean").item(), pyro.param("beta scale").item()))
            Posterior Temperature Mean
            Normal(loc: 19.311052322387695, scale: 0.6258021593093872)
            Posterior Alpha
            Normal(loc: 19.773971557617188, scale: 1.8541947603225708)
            Posterior Beta
```

```
Normal(loc: 1.5178951025009155, scale: 0.1155082955956459)
```

Ice-cream Shop Model:

• We run the (variational) inference engine and get the results.

Pyro

```
M In [48]: #Run inference
           svi(model, guide, obs, num steps=1000)
           #Print results
           print("Posterior Temperature Mean")
           print(dist.Normal(pyro.param("mean").item(), pyro.param("scale").item()))
           print("")
           print("Posterior Alpha")
           print(dist.Normal(pyro.param("alpha mean").item(), pyro.param("alpha scale").item()))
           print("")
           print("Posterior Beta")
           print(dist.Normal(pyro.param("beta mean").item(), pyro.param("beta scale").item()))
            Posterior Temperature Mean
            Normal(loc: 19.311052322387695, scale: 0.6258021593093872)
            Posterior Alpha
            Normal(loc: 19.773971557617188, scale: 1.8541947603225708)
            Posterior Beta
            Normal(loc: 1.5178951025009155, scale: 0.1155082955956459)
```

Ice-cream Shop Model:

- We run the (variational) inference engine and get the results.
- With PPLs, we only care about modeling, **not about the low-level details** of the machine-learning solver.

Exercise 4: Introduce Humidity in the Icecream shop model

- Assume we also have a bunch of humidity senor measurements.
- Assume the sales are also linearly influenced by the humidity.
- Extend the above model in order to integrate all of that. Day1/students_PPLs_Intro.ipynb

Day 2 Probabilistic Models with Deep Neural Networks

Andrés Masegosa

Department of Mathematics University of Almería Spain

Masegosa, A. R., Cabañas, R., Langseth, H., Nielsen, T. D., & Salmern, A. (2019). Probabilistic Models with Deep Neural Networks. arXiv preprint arXiv:1908.03442.

Day 1: Introduction to probabilistic programming languages

- Why do we need PPLs?
- Probabilistic programming in Pyro
- Hand-on exercises:
 - Probability Distributions in Pyro.
 - Probabilistic Models in Pyro.
 - Ice-cream Shop Model.

Day 2: Probabilistic Models with Deep Neural Networks

- Uncertainty in Machine Learning
- Variational Inference
- Supervised/Unsupervised Learning
- Hand-on exercises
 - Bayesian Neural Networks
 - Variational Auto-encoders

Uncertainty with Machine Learning

Uncertainty with Machine Learning



Why is important to model uncertainty?

- To assess confidence in the predictions.
- To know what I don't know.
- Big implications in many problems.



Why is important to model uncertainty?

- To assess confidence in the predictions.
- To know what I don't know.
- Big implications in many problems.

Uncertainty with Machine Learning



Why is important to model uncertainty?

- To assess confidence in the predictions.
- To know what I don't know.
- Big implications in many problems.



Probabilistic Models naturally quantify uncertainty

- Everything is defined in terms of random variables.
- Standard language to model uncertainty.



Probabilistic Models naturally quantify uncertainty

- Everything is defined in terms of random variables.
- Standard language to model uncertainty.

Example

Nunn, N. & Puga, D., Ruggedness: The blessing of bad geography in Africa, Review of Economics and Statistics 94(1), Feb. 2012



Relationship between topographic heterogeneity and GDP per capita

• Is Terrain ruggedness or **bad geography** is related to poorer economic performance outside of Africa?

Nunn, N. & Puga, D., Ruggedness: The blessing of bad geography in Africa, Review of Economics and Statistics 94(1), Feb. 2012



Relationship between topographic heterogeneity and GDP per capita

- Is Terrain ruggedness or **bad geography** is related to poorer economic performance outside of Africa?
- Does rugged terrains have had a reverse effect on income for African nations?

 $y = w_0 + w_1 \cdot x$

$$y = w_0 + w_1 \cdot x$$

Mean Square Error Loss:

$$L(\theta) = \sum_{i=1}^{n} (y_i - w_0 + w_1 \cdot x_i)^2$$

$$y = w_0 + w_1 \cdot x$$

Mean Square Error Loss:

$$L(\theta) = \sum_{i=1}^{n} (y_i - w_0 + w_1 \cdot x_i)^2$$

Image: Second S

$$\arg\min_{\theta} L(\theta)$$
$$\theta^{t+1} = \theta^t - \nabla_{\theta} L(\theta)$$

$$y = w_0 + w_1 \cdot x$$

Mean Square Error Loss:

$$L(\theta) = \sum_{i=1}^{n} (y_i - w_0 + w_1 \cdot x_i)^2$$

Image: Second S

$$\arg\min_{\theta} L(\theta)$$
$$\theta^{t+1} = \theta^t - \nabla_{\theta} L(\theta)$$

Day2/students_Bayesian_regression.ipynb



- Negative slope for Non African Nations.
- Positive slope for African Nations.



- Negative slope for Non African Nations.
- Positive slope for African Nations.

Are we 100% sure about the conclusion?

Bayesian Machine Learning

9 Probabilistic Linear Regression Model:

$$y_i | x_i, \mathbf{w}, \sigma^2 \sim \mathcal{N}(\mu = w_0 + w_1 \cdot x_i, \sigma^2)$$

9 Probabilistic Linear Regression Model:

$$y_i | x_i, \mathbf{w}, \sigma^2 \sim \mathcal{N}(\mu = w_0 + w_1 \cdot x_i, \sigma^2)$$

Ø Model Parameters are Random Variables:

$$w_0, w_1 \sim \mathcal{N}(\mu = 0, \sigma^2 = 100)$$

 $\frac{1}{\sigma^2} \sim \mathcal{G}(\alpha = 1, \beta = 1)$

O Probabilistic Linear Regression Model:

$$y_i | x_i, \mathbf{w}, \sigma^2 \sim \mathcal{N}(\mu = w_0 + w_1 \cdot x_i, \sigma^2)$$

Ø Model Parameters are Random Variables:

$$w_0, w_1 \sim \mathcal{N}(\mu = 0, \sigma^2 = 100)$$

 $\frac{1}{\sigma^2} \sim \mathcal{G}(\alpha = 1, \beta = 1)$

We compute the posterior (Bayes' rule):

$$p(w_0, w_1, \sigma^2 | D) = \frac{\prod_{i=1}^n p(y_i | x_i, w_0, w_1, \sigma^2) p(w_0, w_1, \sigma^2)}{p(y_1, \dots, y_n | x_1, \dots, x_n)}$$

Probabilistic Linear Regression Model:

$$y_i | x_i, \mathbf{w}, \sigma^2 \sim \mathcal{N}(\mu = w_0 + w_1 \cdot x_i, \sigma^2)$$

Ø Model Parameters are Random Variables:

$$w_0, w_1 \sim \mathcal{N}(\mu = 0, \sigma^2 = 100)$$

 $\frac{1}{\sigma^2} \sim \mathcal{G}(\alpha = 1, \beta = 1)$

We compute the posterior (Bayes' rule):

$$p(w_0, w_1, \sigma^2 | D) = \frac{\prod_{i=1}^n p(y_i | x_i, w_0, w_1, \sigma^2) p(w_0, w_1, \sigma^2)}{p(y_1, \dots, y_n | x_1, \dots, x_n)}$$

$Day2/students_Bayesian_regression.ipynb$

Exercise: Bayesian Logistic Regression

• Predicts whether a country is African or not based on ruggedness and GDP.

Day2/students_Bayesian_Regression.ipynb

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

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

• Inference about the unknowns is through the **posterior**, the conditional distribution of the hidden variables given the observations

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}$$

The posterior quantify our uncertainty in the model.

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

• Inference about the unknowns is through the **posterior**, the conditional distribution of the hidden variables given the observations

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}$$

The posterior quantify our uncertainty in the model.

• For most interesting models, the denominator is not tractable.

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

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

• Inference about the unknowns is through the **posterior**, the conditional distribution of the hidden variables given the observations

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}$$

The posterior quantify our uncertainty in the model.

• For most interesting models, the denominator is not tractable.

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z}) p(\mathbf{z}) d\mathbf{z}$$

We have to use approximate posterior inference.

Variational Inference



• Variational Inference is an approximate method,

 $q(\mathbf{z}|\boldsymbol{\nu}) \approx p(\mathbf{z}|\mathbf{x})$

• $q(\mathbf{z}|\boldsymbol{\nu})$ is a simpler distribution (e.g. diagonal covariance matrix, unimodal).



• Variational Inference is an approximate method,

 $q(\mathbf{z}|\boldsymbol{\nu}) \approx p(\mathbf{z}|\mathbf{x})$

- $q(\mathbf{z}|\boldsymbol{\nu})$ is a simpler distribution (e.g. diagonal covariance matrix, unimodal).
- Find the best possible approximation,

$$\boldsymbol{\nu}^{\star} = \arg\min_{\boldsymbol{\nu}} KL(q(\mathbf{z}|\boldsymbol{\nu})||p(\mathbf{z}|\mathbf{x}))$$


• Variational Inference is an approximate method,

 $q(\mathbf{z}|\boldsymbol{\nu}) \approx p(\mathbf{z}|\mathbf{x})$

- $q(\mathbf{z}|\boldsymbol{\nu})$ is a simpler distribution (e.g. diagonal covariance matrix, unimodal).
- Find the best possible approximation,

$$\boldsymbol{\nu}^{\star} = \arg\min_{\boldsymbol{\nu}} KL(q(\mathbf{z}|\boldsymbol{\nu})||p(\mathbf{z}|\mathbf{x}))$$

• Solve using a optimization algorithm (i.e. gradient descent).



$$\boldsymbol{\nu}^{\star} = \arg\min_{\boldsymbol{\nu}} KL(q(\mathbf{z}|\boldsymbol{\nu})||p(\mathbf{z}|\mathbf{x}))$$

$$\boldsymbol{\nu}^{\star} = \arg\min_{\boldsymbol{\nu}} KL(q(\mathbf{z}|\boldsymbol{\nu})||p(\mathbf{z}|\mathbf{x}))$$

$$\boldsymbol{\nu}^{\star} = \arg\min_{\boldsymbol{\nu}} KL(q(\mathbf{z}|\boldsymbol{\nu})||p(\mathbf{z}|\mathbf{x}))$$

• After some manipulations,

$$\ln p(\mathbf{x}) = \mathcal{L}(\boldsymbol{\nu}) - KL(q(\mathbf{z}|\boldsymbol{\nu})||p(\mathbf{z}|\mathbf{x}))$$

$$\boldsymbol{\nu}^{\star} = \arg\min_{\boldsymbol{\nu}} KL(q(\mathbf{z}|\boldsymbol{\nu})||p(\mathbf{z}|\mathbf{x}))$$

• After some manipulations,

$$\ln p(\mathbf{x}) = \mathcal{L}(\boldsymbol{\nu}) - KL(q(\mathbf{z}|\boldsymbol{\nu})||p(\mathbf{z}|\mathbf{x}))$$

• $\mathcal{L}(\boldsymbol{\nu})$ is known as the **ELBO function**.

$$\ln p(\mathbf{x}) \ge \mathcal{L}(\boldsymbol{\nu}) = \mathbb{E}_q[\ln p(\mathbf{x}|\mathbf{z})] - KL(q(\mathbf{z}|\boldsymbol{\nu})||p(\mathbf{z}))$$

$$\boldsymbol{\nu}^{\star} = \arg\min_{\boldsymbol{\nu}} KL(q(\mathbf{z}|\boldsymbol{\nu})||p(\mathbf{z}|\mathbf{x}))$$

• After some manipulations,

$$\ln p(\mathbf{x}) = \mathcal{L}(\boldsymbol{\nu}) - KL(q(\mathbf{z}|\boldsymbol{\nu})||p(\mathbf{z}|\mathbf{x}))$$

• $\mathcal{L}(\boldsymbol{\nu})$ is known as the **ELBO function**.

$$\ln p(\mathbf{x}) \ge \mathcal{L}(\boldsymbol{\nu}) = \mathbb{E}_q[\ln p(\mathbf{x}|\mathbf{z})] - KL(q(\mathbf{z}|\boldsymbol{\nu})||p(\mathbf{z}))$$

• The above minimization problem is equivalent to maximize the ELBO function,

$$\boldsymbol{\nu}^{\star} = \arg \max_{\boldsymbol{\nu}} \mathcal{L}(\boldsymbol{\nu})$$

 $q(\mathbf{z}|\boldsymbol{\nu})$

• $q(\mathbf{z}|\boldsymbol{\nu})$ is a simpler distribution (e.g. diagonal covariance matrix, unimodal).

 $q(\mathbf{z}|\boldsymbol{\nu})$

- $q(\mathbf{z}|\boldsymbol{\nu})$ is a simpler distribution (e.g. diagonal covariance matrix, unimodal).
- **2** Solve this maximization problem,

$$\boldsymbol{\nu}^{\star} = \arg \max_{\boldsymbol{\nu}} \mathcal{L}(\boldsymbol{\nu})$$

 $q(\mathbf{z}|\boldsymbol{\nu})$

- $q(\mathbf{z}|\boldsymbol{\nu})$ is a simpler distribution (e.g. diagonal covariance matrix, unimodal).
- **2** Solve this maximization problem,

$$\boldsymbol{\nu}^{\star} = \arg \max_{\boldsymbol{\nu}} \mathcal{L}(\boldsymbol{\nu})$$

It can be solved by (stochastic) gradient ascent methods:

$$\boldsymbol{\nu}^{(t+1)} = \boldsymbol{\nu}^{(t)} + \alpha \nabla_{\nu} \mathcal{L}(\boldsymbol{\nu}^{(t)})$$

 $q(\mathbf{z}|\boldsymbol{\nu})$

- $q(\mathbf{z}|\boldsymbol{\nu})$ is a simpler distribution (e.g. diagonal covariance matrix, unimodal).
- Solve this maximization problem,

$$\boldsymbol{\nu}^{\star} = \arg \max_{\boldsymbol{\nu}} \mathcal{L}(\boldsymbol{\nu})$$

It can be solved by (stochastic) gradient ascent methods:

$$\boldsymbol{\nu}^{(t+1)} = \boldsymbol{\nu}^{(t)} + \alpha \nabla_{\nu} \mathcal{L}(\boldsymbol{\nu}^{(t)})$$

Use the variational approximation as a proxy,

$$q(\mathbf{z}|\boldsymbol{\nu}) \approx p(\mathbf{z}|\mathbf{x})$$

Variational Inference in Pyro



Pyro's Model

• We have a probabilistic model,

$$p(temp) = \mathcal{N}(15.0, 2.0)$$
$$p(sensor|temp) = \mathcal{N}(sensor, 2.0)$$

```
1 #The guide
2 def guide(obs):
3 a = pyro.param("mean", torch.tensor(0.0))
4 b = pyro.param("scale", torch.tensor(1.), constraint=constraints.positive)
5 temp = pyro.sample('temp', dist.Normal(a, b))
```

Pyro's Guides

- We want to compute the **posterior**, p(temp|sensor = 18.0)
- We define, $q(temp|oldsymbol{
 u}) = \mathcal{N}(oldsymbol{
 u}_a,oldsymbol{
 u}_b)$

```
1 #The guide
2 def guide(obs):
3 a = pyro.param("mean", torch.tensor(0.0))
4 b = pyro.param("scale", torch.tensor(1.), constraint=constraints.positive)
5 temp = pyro.sample('temp', dist.Normal(a, b))
```

Pyro's Guides

- We want to compute the **posterior**, p(temp|sensor = 18.0)
- We define, $q(temp|oldsymbol{
 u}) = \mathcal{N}(oldsymbol{
 u}_a,oldsymbol{
 u}_b)$
- Guide's requirements:
 - Model and Guide have both same input signature.
 - All unobserved sample statements in the model must appear in the guide.

Pyro's Variational Inference

• We run optimization to solve,

$$\boldsymbol{\nu}^{\star} = \arg\min_{\boldsymbol{\nu}} KL(q(temp|\boldsymbol{\nu})||p(temp|sensor = 18.0))$$

Pyro's Variational Inference

• We run optimization to solve,

$$\boldsymbol{\nu}^{\star} = \arg\min_{\boldsymbol{\nu}} KL(q(temp|\boldsymbol{\nu})||p(temp|sensor = 18.0))$$

• It can be solved by (stochastic) gradient ascent methods:

$$\boldsymbol{\nu}^{(t+1)} = \boldsymbol{\nu}^{(t)} + \alpha \nabla_{\boldsymbol{\nu}} \mathcal{L}(\boldsymbol{\nu}^{(t)})$$

Exercise: Pyro implementation for a simple Gaussian model

• Implement a Pyro's **model** for a Normal distribution with prior over mean and (inverse) of variance. The variable *x* is always observed.

$$\begin{split} p(\mu) &= \mathsf{Normal}(0, 10000.0) \\ p(\gamma) &= \mathsf{Gamma}(1, 1) \\ p(x \,|\, \mu, \gamma) &= \mathsf{Norma}(\mu, 1/\gamma) \end{split}$$

Obfine a Pyro's guide for approximating the posterior,

$$q(\mu) = \operatorname{Norma}(\boldsymbol{\nu}_{\mu}, \boldsymbol{\nu}_{\sigma^2})$$
$$q(\gamma) = \operatorname{Gammal}(\boldsymbol{\nu}_{\alpha}, \boldsymbol{\nu}_{\beta})$$



Neural Network

• $h(\cdot; \theta)$ encodes a neural network (i.e. a non-linear function).

 $\mathbf{y} = h(\mathbf{x}; \theta)$

Loss minimization,

$$\arg\min_{\theta} \sum_{i=1}^{n} \ell(h(\mathbf{x}_i; \theta), \mathbf{y}_i)$$



Bayesian Neural Network

• E.g. Neural regressor,

$$p(y|\mathbf{x}, \theta) = \mathcal{N}(\mu = h(\mathbf{x}; \theta_h), \sigma^2 = g(\mathbf{x}; \theta_g))$$

 $h(\cdot; \theta)$ encodes a neural network.

Supervised Learning



Bayesian Neural Network

Supervised Learning



Bayesian Neural Network

• Variational Inference,

$$\boldsymbol{\nu}^{\star} = \arg\min_{\boldsymbol{\nu}} KL(q(\boldsymbol{\theta}|\boldsymbol{\nu})||p(\boldsymbol{\theta}|\mathbf{y}_{1},\mathbf{x}_{1},\ldots,\mathbf{y}_{n},\mathbf{x}_{n}))$$



Bayesian Neural Network

• Variational Inference,

$$\boldsymbol{\nu}^{\star} = \arg\min_{\boldsymbol{\nu}} KL(q(\boldsymbol{\theta}|\boldsymbol{\nu})||p(\boldsymbol{\theta}|\mathbf{y}_{1},\mathbf{x}_{1},\ldots,\mathbf{y}_{n},\mathbf{x}_{n}))$$

• Bayesian predictive is an ensemble of neural networks,

$$p(y|\mathbf{x}) = \int p(y|\mathbf{x}, \theta) q(\theta|\boldsymbol{\nu}^{\star}) d\theta$$
$$\approx \frac{1}{M} \sum_{\theta_j \sim q(\theta|\boldsymbol{\nu}^{\star})} p(y|\mathbf{x}, \theta_j)$$



Bayesian Neural Networks

- Capture the uncertainty.
- Better identify low confidence predictions (They know what they don't know).



Bayesian Neural Networks

- Capture the uncertainty.
- Better identify low confidence predictions (They know what they don't know).

Day2/BayesianNeuralNetworks.ipynb

Unsupervised Probabilistic Models with Deep Neural Networks



Unsepervised Learning with Deep Neural Networks

• No labeled data (extract structure from data).



Unsepervised Learning with Deep Neural Networks

- No labeled data (extract structure from data).
- A generative model,

 $\mathbf{z} \sim p(\mathbf{z})$ $\mathbf{x} \sim p(\mathbf{x}|\mathbf{z}, \theta)$



Unsepervised Learning with Deep Neural Networks

- No labeled data (extract structure from data).
- A generative model,

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

 $\mathbf{x} \sim p(\mathbf{x}|\mathbf{z}, \theta)$

• Variational Inference,

 $q(\mathbf{z}|\mathbf{x}, \phi) \approx p(\mathbf{z}|\mathbf{x}, \theta)$

Unsupervised Learning



Unsepervised Learning with Deep Neural Networks

- No labeled data (extract structure from data).
- A generative model,

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

 $\mathbf{x} \sim p(\mathbf{x}|\mathbf{z}, \theta)$

• Variational Inference,

$$q(\mathbf{z}|\mathbf{x}, \phi) \approx p(\mathbf{z}|\mathbf{x}, \theta)$$



Kingma, D. P., & Welling, M. (2019). An Introduction to Variational Autoencoders. arXiv preprint arXiv:1906.02691.



(White et al. 2016)

Exercise: Bayesian PCA in Pyro

- VAE using linear transformations.
- Applied to MNIST data set.

Day2/BayesianPCA.ipynb

Exercise: Bayesian PCA in Pyro

- VAE using linear transformations.
- Applied to MNIST data set.

Day2/BayesianPCA.ipynb

Exercise: Variational Auto-Encoder in Pyro

- VAE using simple neural network.
- Applied to MNIST data set.

Day2/VAE.ipynb

• Probabilistic modeling is a key aspect in AI systems

- Confidence in model predictions/actions.
- To know what you do not know.

• Probabilistic modeling is a key aspect in AI systems

- Confidence in model predictions/actions.
- To know what you do not know.

• PPLs are the right tool for probabilistic modeling.

- Enormous expressibility.
- Powerful inference engines (BlackBox Variational Inference).



https://inferpy.readthedocs.io



http://www.amidsttoolbox.com/

Thanks!!!!! :)