# Bayesian ML

## A practical approach

**Verena, Filip**

# Motivation

Key issue of ML: **<u>overfitting</u>**

→ solution at hand: regularisation (dropout,…)

Became concern for us with **<u>Model-based Reinforcement Learning</u>**.

Disclaimer: only in the research phase. No successful application for RL yet.

Filip will show you another application.

# Overfitting

Example regression:

NN is a probabilistic model $p(y\,|\,\mathbf{x}, \mathbf{w})$, $p$ is a Gaussian distribution.

With training data $\mathscr{D} = \{x^{(i)}, y^{(i)}\}$ optimise $\mathbf{w}$ by maximising likelihood function $p(\mathscr{D}\,|\,\mathbf{w}) = \prod_i p(y^{(i)}\,|\,\mathbf{x}^{(i)}, \mathbf{w})$... Maximum

Likelihood Estimate $\rightarrow$ <span style="color:red">point estimate of $\mathbf{w}$</span>

Better: regularisation

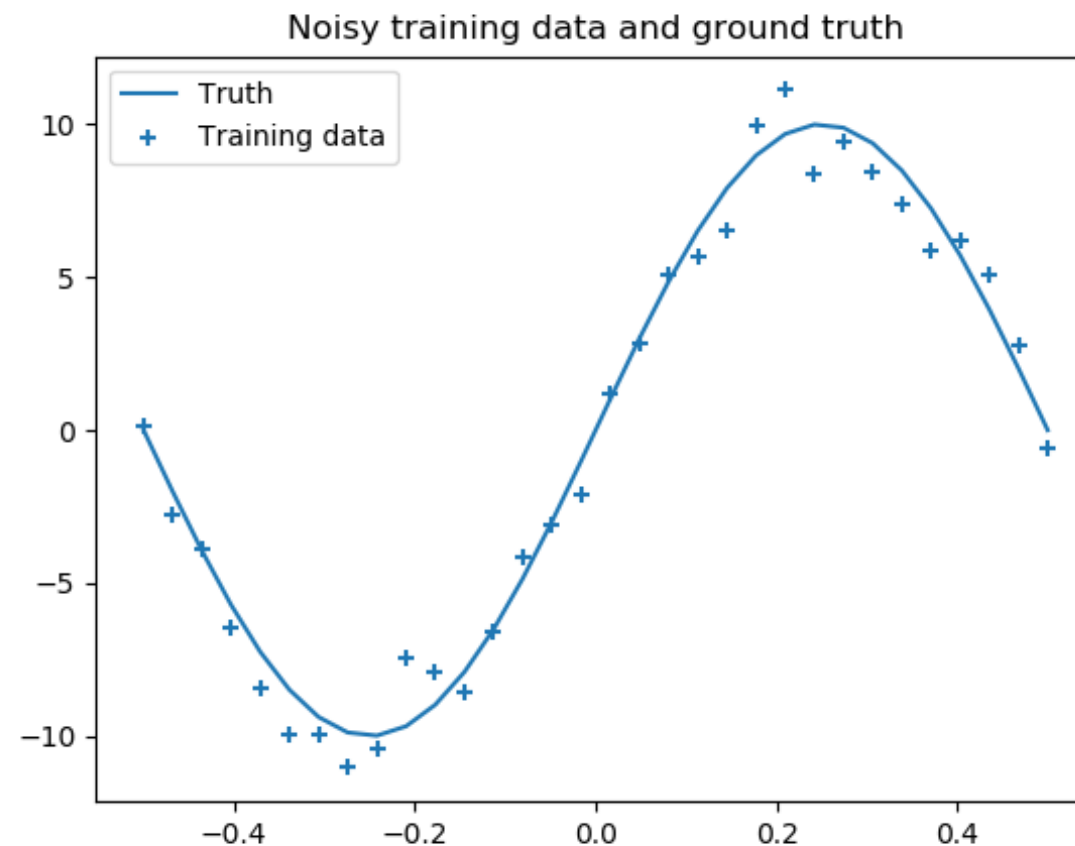Would like to maximise $p(\mathbf{w}\,|\,\mathscr{D})$ ...maximum a posteriori estimate of $\mathbf{w}$

Do this with $p(\mathbf{w}\,|\,\mathscr{D}) \propto p(\mathscr{D}\,|\,\mathbf{w})p(\mathbf{w})$... prior distribution $p(\mathbf{w})$.

<span style="color:red">$\rightarrow$ still only point estimate of $\mathbf{w}$</span>

# Example

Dataset is 32 points of noisy sinusoidal function.

```python
18  def f(x, sigma):
19      epsilon = np.random.randn(*x.shape) * sigma
20      return 10 * np.sin(2 * np.pi * (x)) + epsilon
21
22  train_size = 32
23  noise = 1.0
24
25  X = np.linspace(-0.5, 0.5, train_size).reshape(-1, 1)
26  y = f(X, sigma=noise)
27  y_true = f(X, sigma=0.0)
```



Noisy training data and ground truth

# Example

Simple feedforward regression model with regularisation with TensorFlow 2

```
1    import tensorflow.keras as tfk
2    from tensorflow.keras import initializers, activations
3    from tensorflow.keras.layers import Layer
4    import tensorflow_probability as tfp
5    import tensorflow as tf
6    tfd = tfp.distributions
7
```
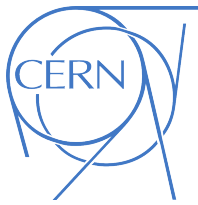
```
41   batch_size = train_size
42   num_batches = train_size / batch_size
43
44
45
46   # Build model.
47   model = tf.keras.Sequential([
48       tfk.layers.Input(shape=(1,)),
49       tfk.layers.Dense(20, activation="relu"),
50       tfk.layers.Dropout(0.1),
51       tfk.layers.Dense(20, activation="relu"),
52       tfk.layers.Dropout(0.1),
53   tfk.layers.Dense(1)
54   ,tfk.layers.Dropout(0.1)
55   ])
56
57
58
59   from tensorflow.keras import callbacks, optimizers
60
61   def neg_log_likelihood(y_obs, y_pred, sigma=noise):
62       dist = tfp.distributions.Normal(loc=y_pred, scale=sigma)
63       return tfk.backend.sum(-dist.log_prob(y_obs))
64
65   model.compile(loss=neg_log_likelihood, optimizer=optimizers.Adam(lr=0.08), metrics=['mse'])
66   model.fit(X, y, batch_size=batch_size, epochs=1500, verbose=1);
67
```

**Model**

⟵ **Assume** $p(y\,|\,\mathbf{x},\mathbf{w})$ **is Gaussian**

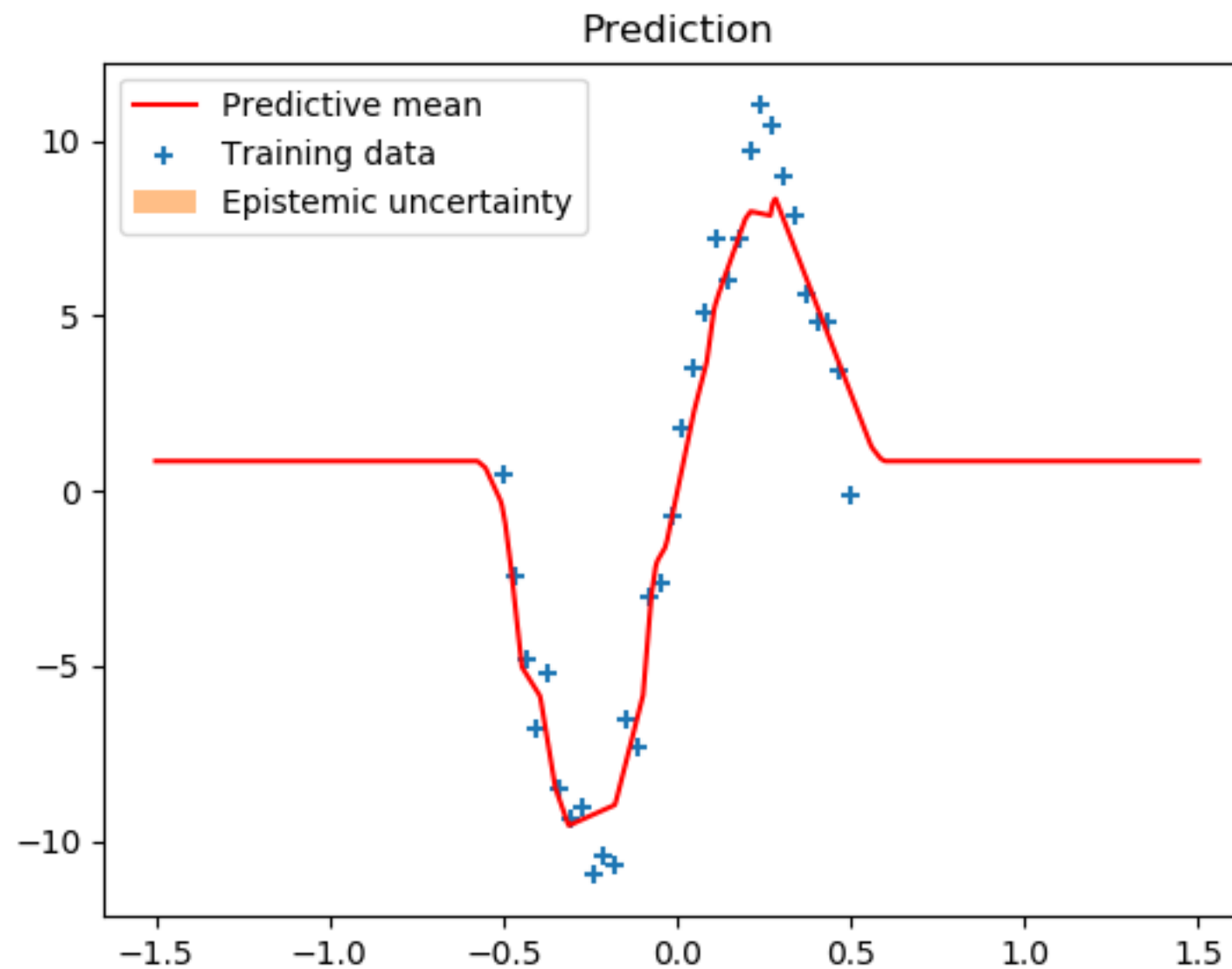ML coffee, 8/5/2020

# Example

## Test of model:

```
68
69      X_test = np.linspace(-1.5, 1.5, 1000).reshape(-1, 1)
70      y_pred_list = []
71
72      import tqdm
73
74      for i in tqdm.tqdm(range(500)):
75          y_pred = model.predict(X_test)
76          y_pred_list.append(y_pred)
77
78      y_preds = np.concatenate(y_pred_list, axis=1)
79
80      y_mean = np.mean(y_preds, axis=1)
81      y_sigma = np.std(y_preds, axis=1)
82
83      plt.plot(X_test, y_mean, 'r-', label='Predictive mean');
84      plt.scatter(X, y, marker='+', label='Training data')
85      plt.fill_between(X_test.ravel(),
86                       y_mean + 2 * y_sigma,
87                       y_mean - 2 * y_sigma,
88                       alpha=0.5, label='Epistemic uncertainty')
89      plt.title('Prediction')
90      plt.legend();
91
92      plt.show()
```

# Example

Result:



Can we trust this? Need **epistemic uncertainty** on predictions.

# Towards Variational Layers

With regularisation already assume distributions of $\mathbf{w}$…how to extract uncertainty on predictions of NN?

$\rightarrow$ Average predictions over ensemble of NNs weighted by posterior probabilities of their $\mathbf{w}$

$$p(y \,|\, \mathbf{x}, \mathcal{D}) = \int p(y \,|\, \mathbf{x}, \mathbf{w}) p(\mathbf{w} \,|\, \mathcal{D}) d\mathbf{w}$$

But we cannot write $p(\mathbf{w} \,|\, \mathcal{D})$ down.

$\rightarrow$ Use variational distribution $q(\mathbf{w} \,|\, \theta)$ of known functional form. Need to find $\theta$ by minimising KL divergence between $p(\mathbf{w} \,|\, \mathcal{D})$ and $q(\mathbf{w} \,|\, \theta)$
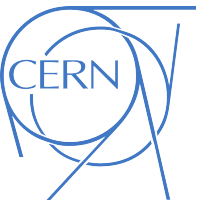
# Towards Variational Layers

Resulting cost function of KL divergence minimisation can be arranged as:

$$\mathscr{F}(\mathscr{D}, \theta) = \mathbb{E}_{q(\mathbf{w}|\theta)} \log q(\mathbf{w}|\theta) - \mathbb{E}_{q(\mathbf{w}|\theta)} \log p(\mathbf{w}) - \mathbb{E}_{q(\mathbf{w}|\theta)} \log p(\mathscr{D}|\mathbf{w})$$

**Data independent**

See ELBO in *https://medium.com/tensorflow/ regression-with-probabilistic-layers-in-tensorflow-probability-e46ff5d37baf*

# Constructing a Variational Layer yourself

Assume Gaussian as variational posterior for $\mathbf{w}$
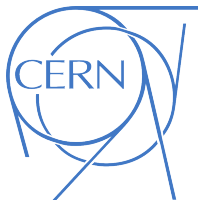
$$\rightarrow \theta = (\mu, \sigma)$$

Assume non-trainable prior:
$$p(\mathbf{w}) = \pi \mathcal{N}(\mathbf{w}\,|\,0,\sigma_1) + (1 - \pi)\mathcal{N}(\mathbf{w}\,|\,0,\sigma_2)$$

Basic steps:

- define your trainable parameters in `build()`

- define in `call()` how to evaluate your layer given `input`

# Constructing a Variational Layer yourself

Example:

```python
class MyDenseVariational(Layer):
    def __init__(self,
                 units,
                 kl_weight,
                 activation=None,
                 prior_sigma_1=1.5,
                 prior_sigma_2=0.1,
                 prior_pi=0.5, **kwargs):
        self.units = units
        self.kl_weight = kl_weight
        self.activation = activations.get(activation)
        self.prior_sigma_1 = prior_sigma_1
        self.prior_sigma_2 = prior_sigma_2
        self.prior_pi_1 = prior_pi
        self.prior_pi_2 = 1.0 - prior_pi
        self.init_sigma = np.sqrt(self.prior_pi_1 * self.prior_sigma_1 ** 2 +
                                  self.prior_pi_2 * self.prior_sigma_2 ** 2)

        super().__init__(**kwargs)

    def compute_output_shape(self, input_shape):...

    def build(self, input_shape):...

    def call(self, inputs, **kwargs):...

    def kl_loss(self, w, mu, sigma):...

    def log_prior_prob(self, w):...
```
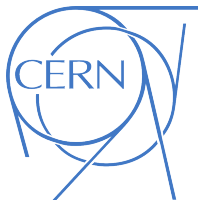
# Constructing a Variational Layer yourself

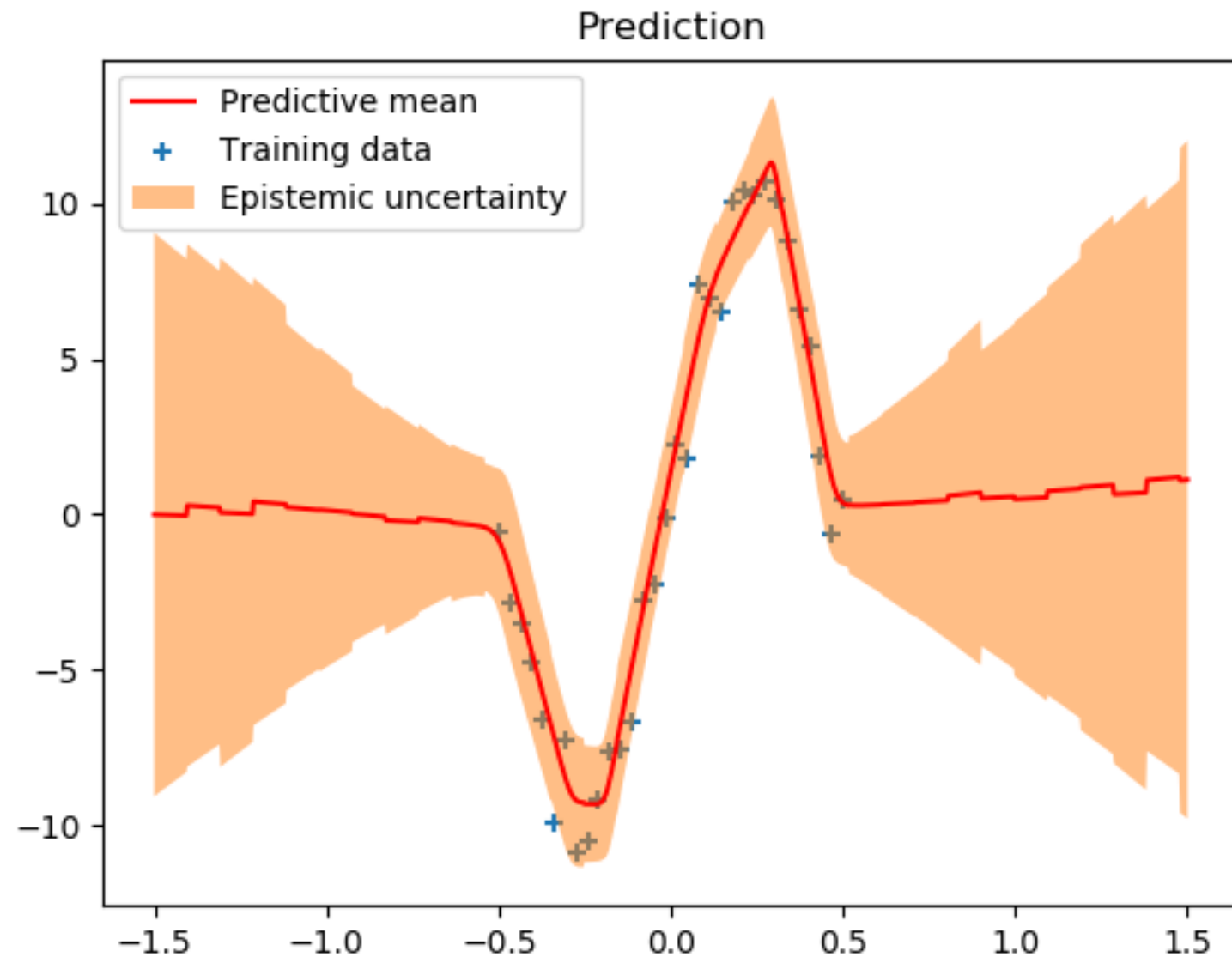For the details of `build()` and `call()`:

```python
36      def build(self, input_shape):
37          self.kernel_mu = self.add_weight(name='kernel_mu',
38                                  shape=(input_shape[1], self.units),
39                                  initializer=initializers.RandomNormal(stddev=self.init_sigma),
40                                  trainable=True)
41          self.bias_mu = self.add_weight(name='bias_mu',
42                                  shape=(self.units,),
43                                  initializer=initializers.RandomNormal(stddev=self.init_sigma),
44                                  trainable=True)
45          self.kernel_rho = self.add_weight(name='kernel_rho',
46                                  shape=(input_shape[1], self.units),
47                                  initializer=initializers.constant(0.0),
48                                  trainable=True)
49          self.bias_rho = self.add_weight(name='bias_rho',
50                                  shape=(self.units,),
51                                  initializer=initializers.constant(0.0),
52                                  trainable=True)
53          super().build(input_shape)
54
55      def call(self, inputs, **kwargs):
56          kernel_sigma = tf.math.softplus(self.kernel_rho)
57          kernel = self.kernel_mu + kernel_sigma * tf.random.normal(self.kernel_mu.shape)
58
59          bias_sigma = tf.math.softplus(self.bias_rho)
60          bias = self.bias_mu + bias_sigma * tf.random.normal(self.bias_mu.shape)
61
62          self.add_loss(self.kl_loss(kernel, self.kernel_mu, kernel_sigma) +
63                        self.kl_loss(bias, self.bias_mu, bias_sigma))
64
65          return self.activation(tfk.backend.dot(inputs, kernel) + bias)
66
67      def kl_loss(self, w, mu, sigma):
68          variational_dist = tfp.distributions.Normal(mu, sigma)
69          return self.kl_weight * tfk.backend.sum(variational_dist.log_prob(w) - self.log_prior_prob(w))
70
```
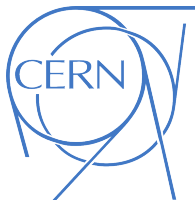
# Using MyDenseVariational

```
105
106    kl_weight = 1.0 / num_batches
107    prior_params = {
108        'prior_sigma_1': 1.5,
109        'prior_sigma_2': 0.1,
110        'prior_pi': 0.5
111    }
112
113    x_in = Input(shape=(1,))
114    x = MyDenseVariational(20, kl_weight, **prior_params, activation='relu')(x_in)
115    x = MyDenseVariational(20, kl_weight, **prior_params, activation='relu')(x)
116    x = MyDenseVariational(1, kl_weight, **prior_params)(x)
117
118    model = Model(x_in, x)
119
```



Prediction

# Using TF2 DenseVariational Layer

```
70      # Build model.
71      model = tf.keras.Sequential([
72          tf.keras.layers.Input(shape=(1,)),
73          tfp.layers.DenseVariational(units=20,
74                                      make_posterior_fn=posterior_mean_field,
75                                      make_prior_fn=prior_trainable,
76                                      kl_weight=kl_weight,
77                                      activation='relu'),
78          tfp.layers.DenseVariational(units=20,
79                                      make_posterior_fn=posterior_mean_field,
80                                      make_prior_fn=prior_trainable,
81                                      kl_weight=kl_weight,
82                                      activation='relu'),
83          tfp.layers.DenseVariational(units=1,
84                                      make_posterior_fn=posterior_mean_field,
85                                      make_prior_fn=prior_trainable,
86                                      kl_weight=kl_weight)
87      ])
88
```

```
# Specify the surrogate posterior over `keras.layers.Dense` `kernel` and `bias`.
def posterior_mean_field(kernel_size, bias_size=0, dtype=None):
    n = kernel_size + bias_size
    c = np.log(np.expm1(1.))
    return tf.keras.Sequential([
        tfp.layers.VariableLayer(2 * n, dtype=dtype),
        tfp.layers.DistributionLambda(lambda t: tfd.Independent(
            tfd.Normal(loc=t[..., :n],
                       scale=1e-5 + tf.nn.softplus(c+t[..., n:])),
            reinterpreted_batch_ndims=1)),
    ])
```
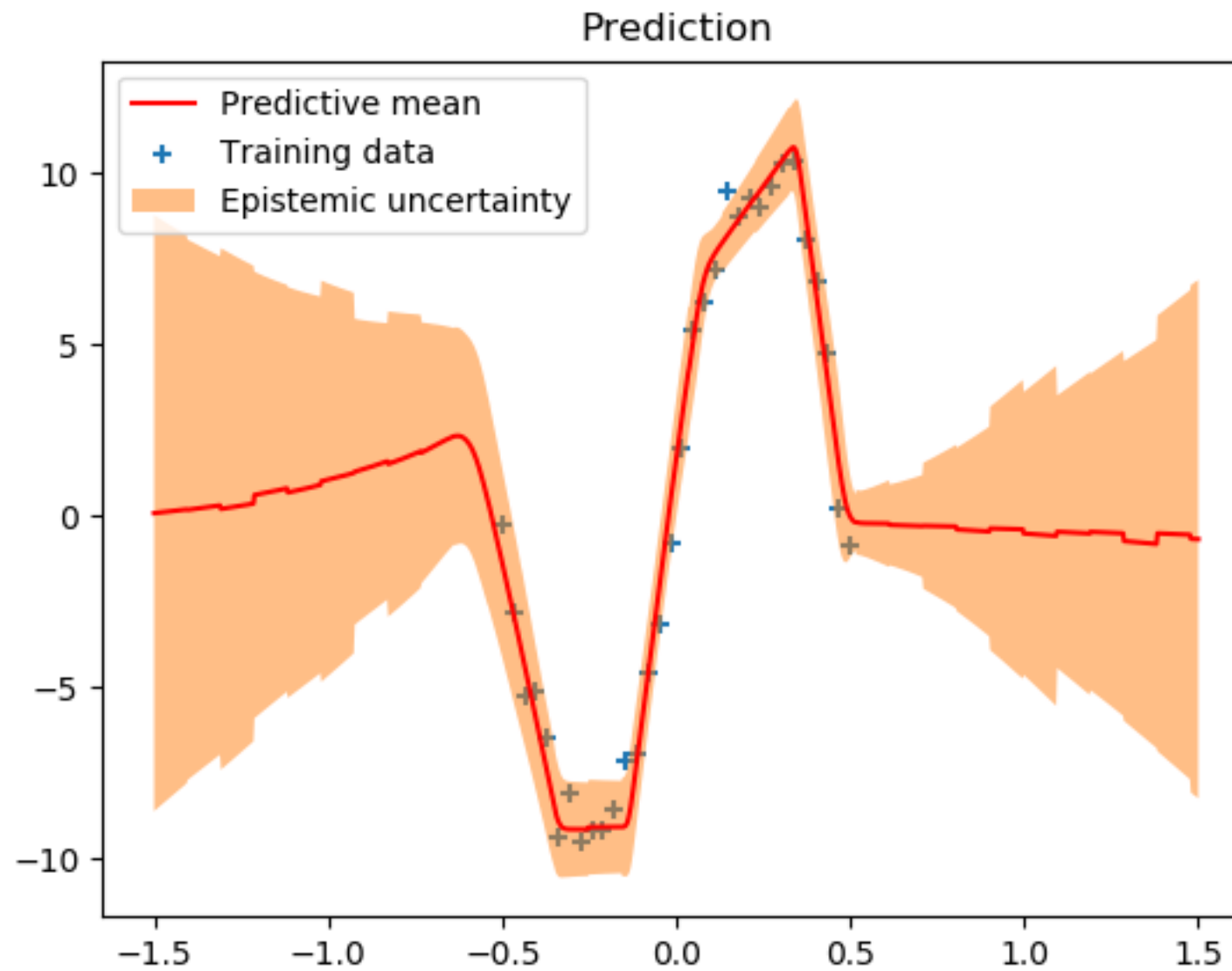
```
# Specify the prior over `keras.layers.Dense` `kernel` and `bias`.
def prior_trainable(kernel_size, bias_size=0, dtype=None):
    n = kernel_size + bias_size
    return tf.keras.Sequential([
        tfp.layers.VariableLayer(n, dtype=dtype),
        tfp.layers.DistributionLambda(lambda t: tfd.Independent(
            tfd.Normal(loc=t, scale=1),
            reinterpreted_batch_ndims=1)),
    ])
```

# Using TF2 DenseVariational Layer
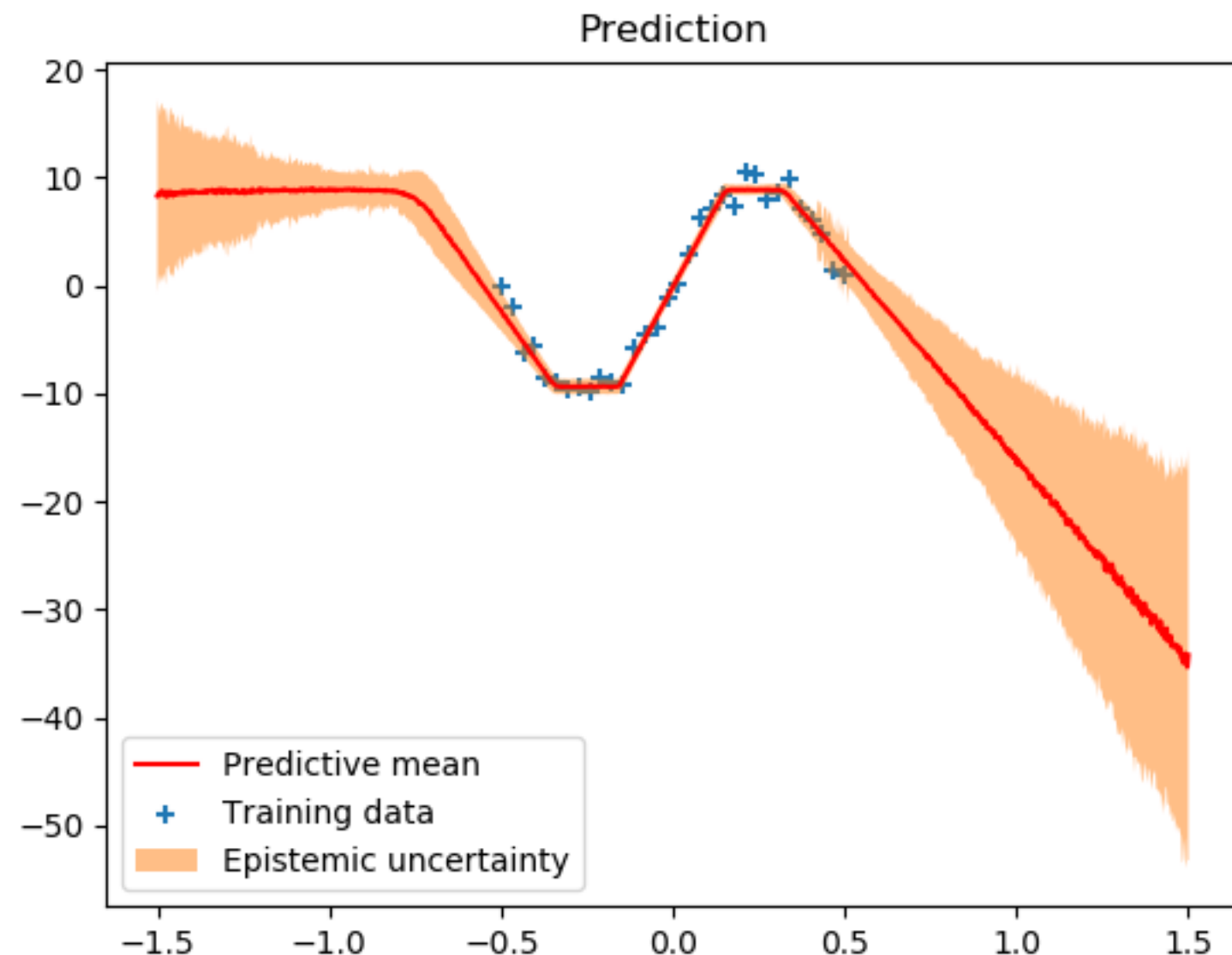
Result:

# DenseFlipout

Assume distribution of weights and biases.

→ prior and posterior distributions for those to be defined.

Use MonteCarlo approach for sampling weights and biases and integrate over it.

```
# Build model.
model = tf.keras.Sequential([
    tf.keras.layers.Input(shape=(1,)),
    tfp.layers.DenseFlipout(20, activation="relu"),
    tfp.layers.DenseFlipout(20, activation="relu"),
tfp.layers.DenseFlipout(1)

])
```

# Further reading...

https://brendanhasz.github.io/2019/07/23/bayesian-density-net.html

https://medium.com/tensorflow/regression-with-probabilistic-layers-in-tensorflow-probability-e46ff5d37baf

https://github.com/krasserm/bayesian-machine-learning