

Uncertainty Quantification in Deep Neural Networks

Alexey Zaytsev

14 March, 2025

Content

- 1. Why do we need neural networks?**
- 2. What is uncertainty**
 - Definition of uncertainty
 - Definition of uncertainty quantification
 - Taxonomy: model and data uncertainty
- 3. Road to uncertainty estimates for Neural Networks:**
 - Linear models
 - Pseudo-Bayesian methods: ensembles and sampling
 - Deterministic uncertainty for neural networks
- 4. Open problems**

Why do we need neural networks?

Classic machine learning, representations are available

x – an object

A human

y – a true label

Will buy a product in 3 months?

$h = g(x)$ – an object representation

Salary, age

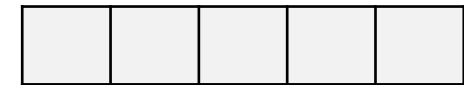
Problem: train a model that can identify the true class label $\approx y$

A logistic regression model $f(h^T \theta)$ outputs a probability of a purchase

Object



Representation



$f(h^T \theta)$

Classifier



For structured data we need representation learning

x – an object

An image

y – a true label

Cat or dog?

$h = g(x)$ – an object representation

NOT CLEAR, HOW TO DEFINE???

Problems: (1) train a model's body that produces representation, (2) train a model's head that predicts the true class;

A neural network!

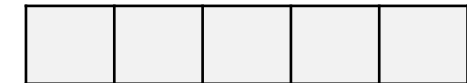
Object



(1)



Representation



(2)



Classifier



Neural Large Language Models

x – an object

A text

y – a true label

The next word

$h = g(x)$ – an object representation

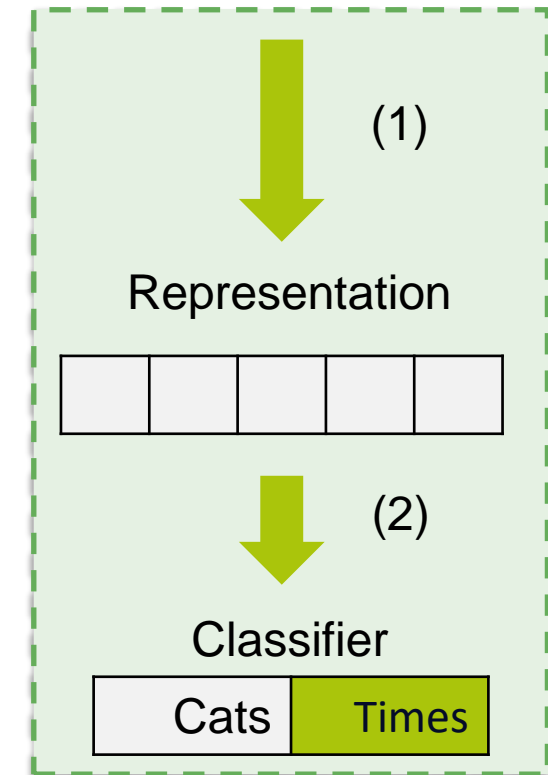
NEED TO TRAIN

Problems: (1) train a model's body that produces representation, (2) train a model's head that predicts the true class;

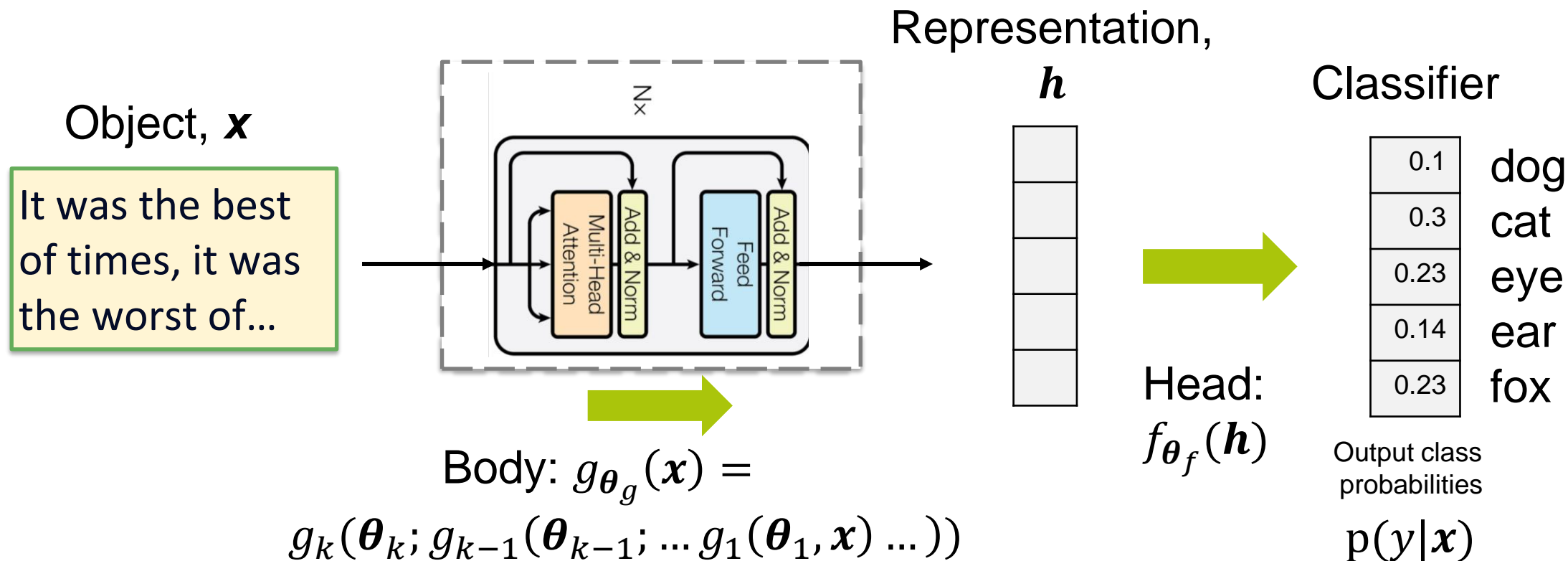
A neural network language model!

Object

It was the best of times, it was the worst of...



Structure of a neural network

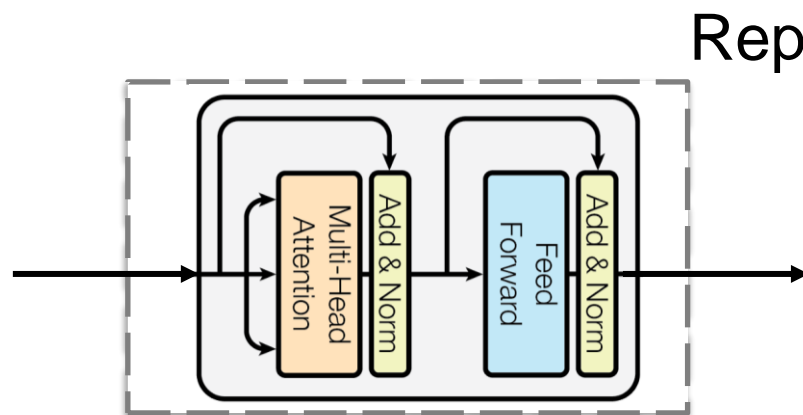


BERT model: 345 millions of parameters in the body
 Largest LLAMA model: 70 billions of parameters

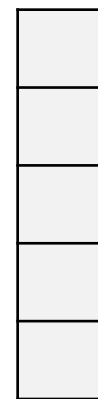
Learning of a neural network

Object, \mathbf{x}

It was the best
of times, it was
the worst of...



Representation, \mathbf{h}



$$f_{\theta_f}(\mathbf{h})$$

Classifier

0.1
0.3
0.23
0.14
0.23

dog
cat
eye
ear
fox

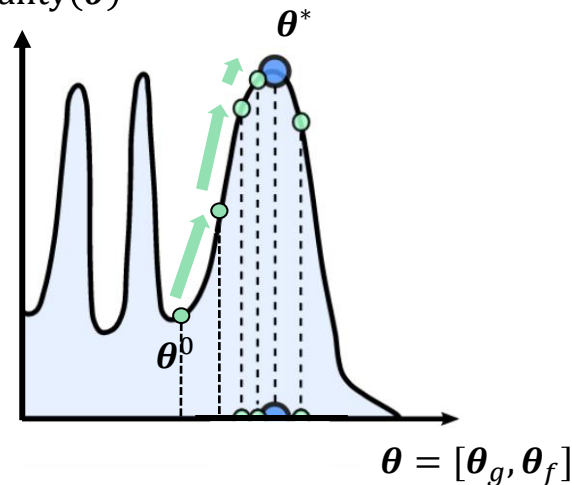
Ground
truth

0
1
0
0
0

$$p(\mathbf{y}|\mathbf{x}) = f_{\theta_f}(\mathbf{h})$$

$$p_{\text{true}}(\mathbf{y}|\mathbf{x})$$

Quality(θ)

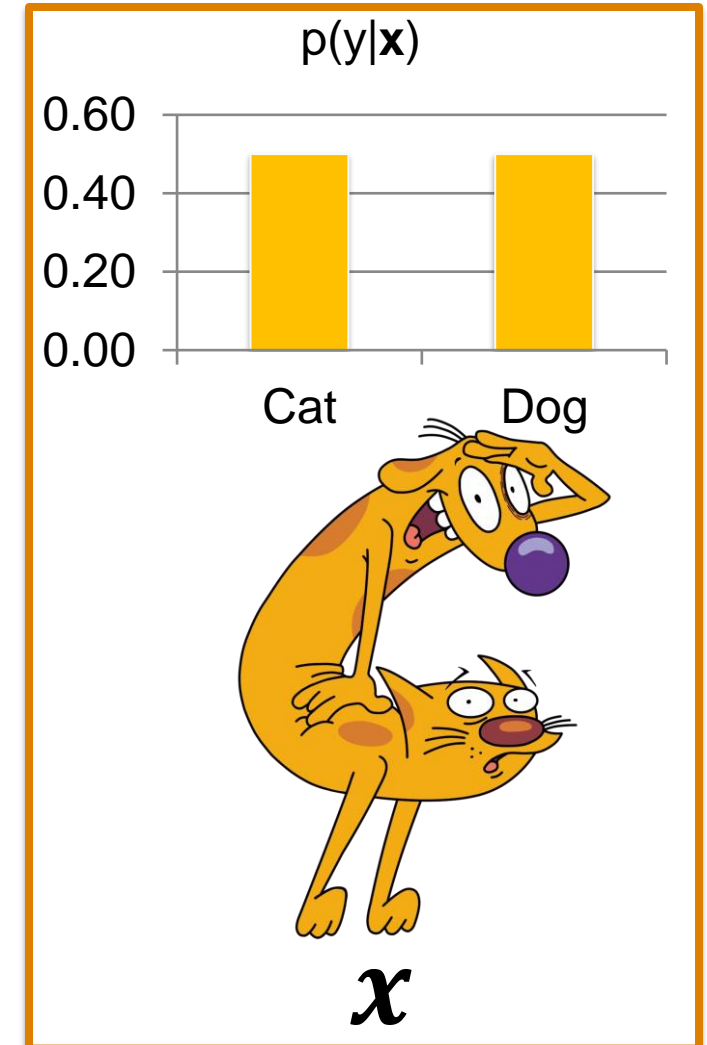
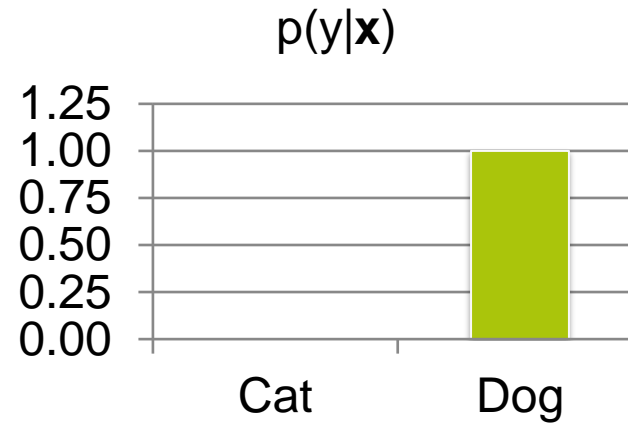
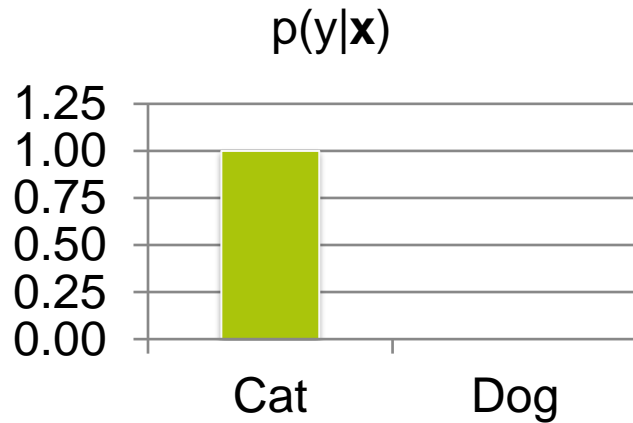


$$g_{\theta_g}(\mathbf{x})$$



Learning is an adjustment of parameters, so the output distribution matches the true distribution

We (almost) never know the uncertainties in the ground truth



Credit: <https://www.jamiesale-cartoonist.com/free-cartoon-dog-vector-clip-art>
<https://www.vecteezy.com/png/13078569-illustration-of-cute-colored-cat-cartoon-cat-image-in-png-format-suitable-for-children-s-book-design-elements-introduction-of-cats-to-children-books-or-posters-about-animal>

What is uncertainty?

What is uncertainty quantification?

Life is uncertain



Uncertainty in daily life examples

- **Weather forecasting** - is it going to rain?
- **Investment decisions** - will the stock go up or down?
- **Medical diagnosis** - what's the exact disease?

What is Uncertainty?

- **Uncertainty** refers to situations involving imperfect or unknown information.
- It's a gap between what we know and what we don't know.

Go back to Neural Networks

- Consider a neural network making predictions: is it a cat or a dog in the image?
- Just like the examples mentioned, the network operates under unknown factors.
- For example, here the network predicts for irrelevant input.

$$P(y = 1|x) = 0.9$$



$$y_{true} = 1$$

$$P(y = 1|x) = 0.2$$



$$y_{true} = 0$$

$$P(y = 1|x) = 0.8$$



$$y_{true} = ???$$

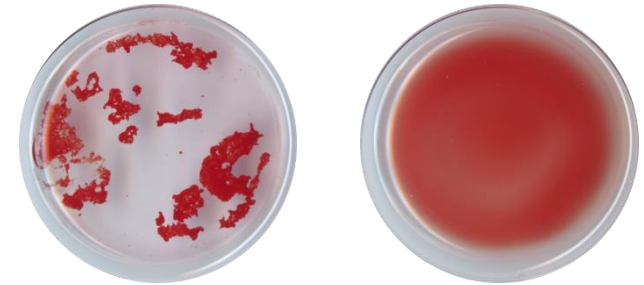
Uncertainty in applications

- **Hallucinations in Large Language models** – how sure is a model about its' output? What hallucinations we fight?
- **Medical diagnosis** – should we retake a performed test?
- **Active learning** – what data should we label to minimize labeling costs, if the labeling budget is limited?

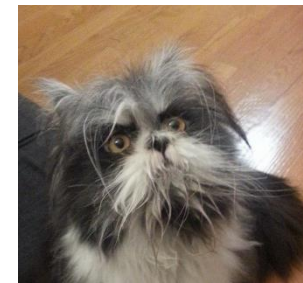
Skoltech



Source: OpenAI



Source: Synteco, IITP RAS



Source:
<https://twitter.com/1evilidiot/status/794613309613821952>

Why Uncertainty Quantification?

- By **quantifying** uncertainty, we can provide a measure of confidence in the network's predictions.
- This measure helps users and developers trust the AI system and make informed decisions based on its outputs.

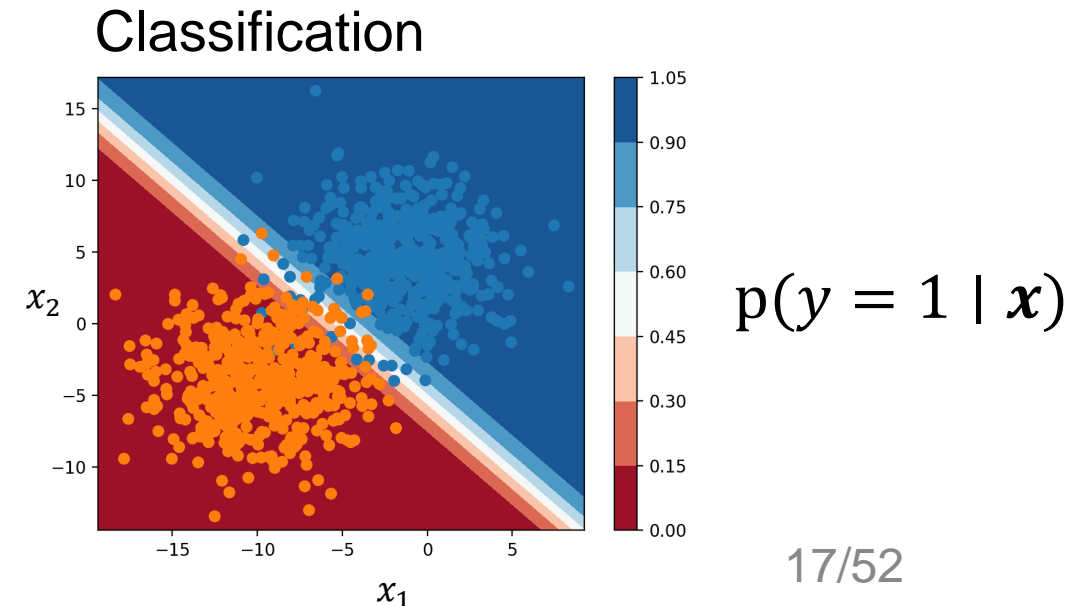
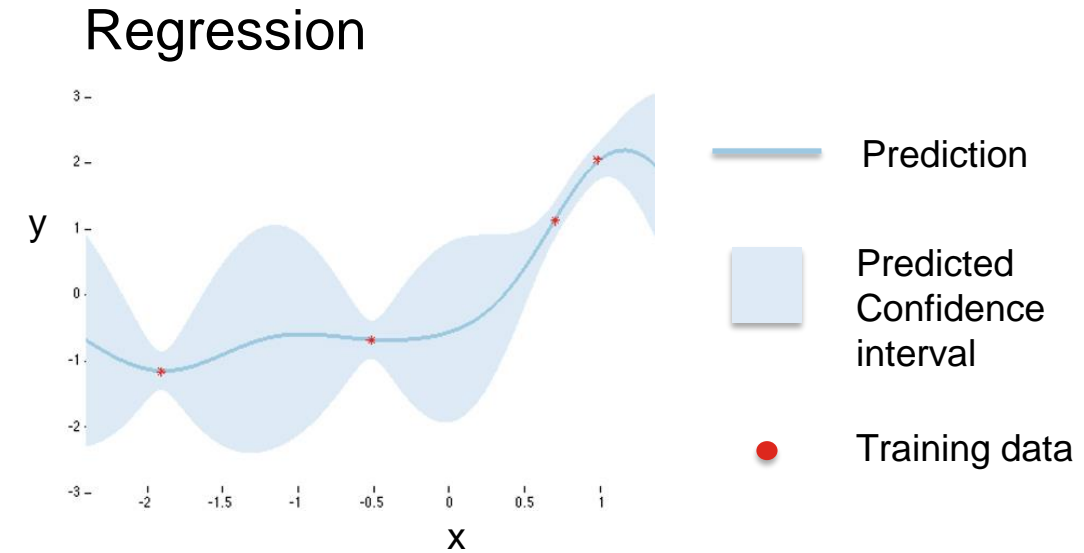
What do we mean by uncertainty?

We would like to have **some value** *distribution over plausible predictions*, or *aggregated value*, which describes the “variability” over possible predictions.

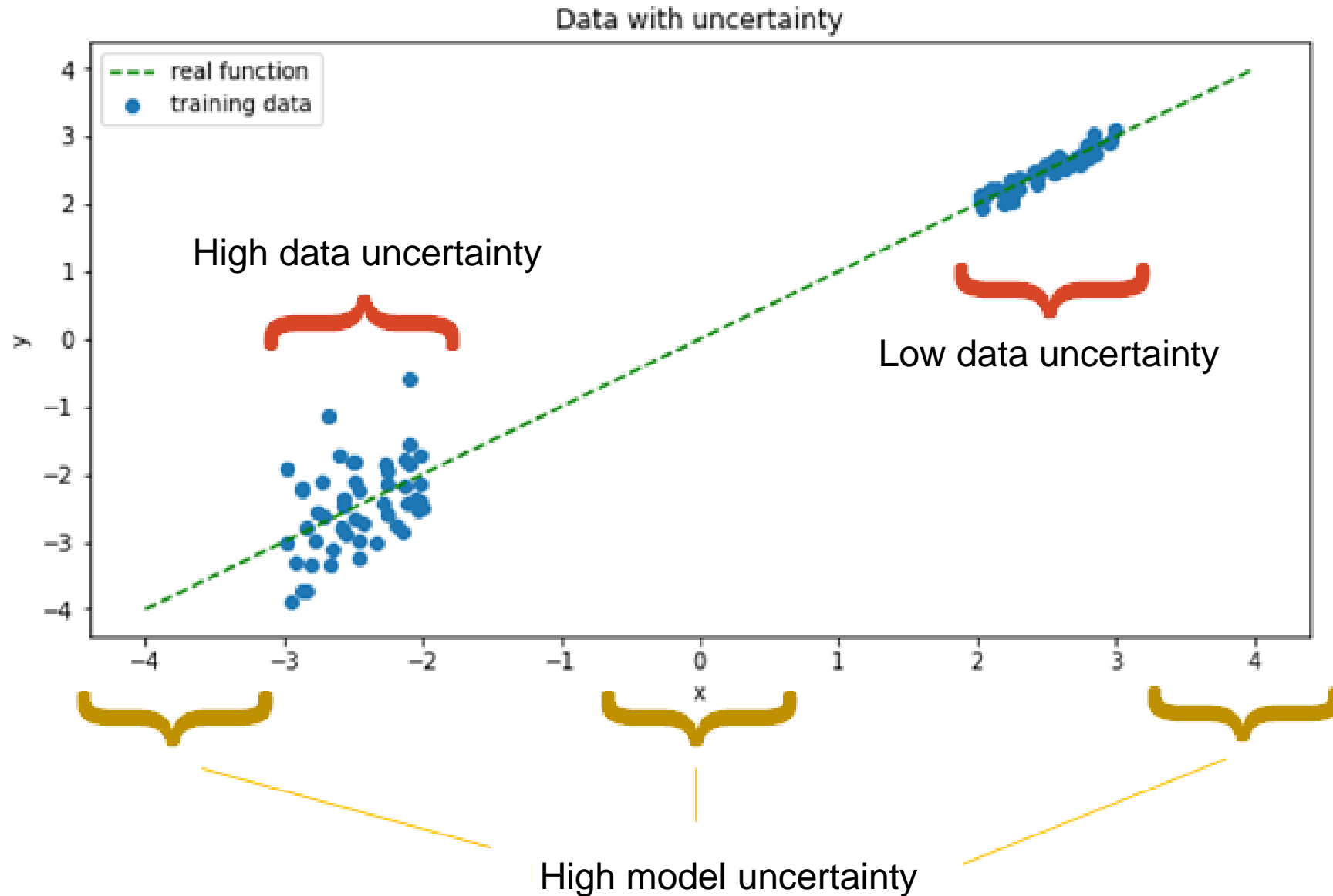
Example:

Regression: predict a mean with *the variance*

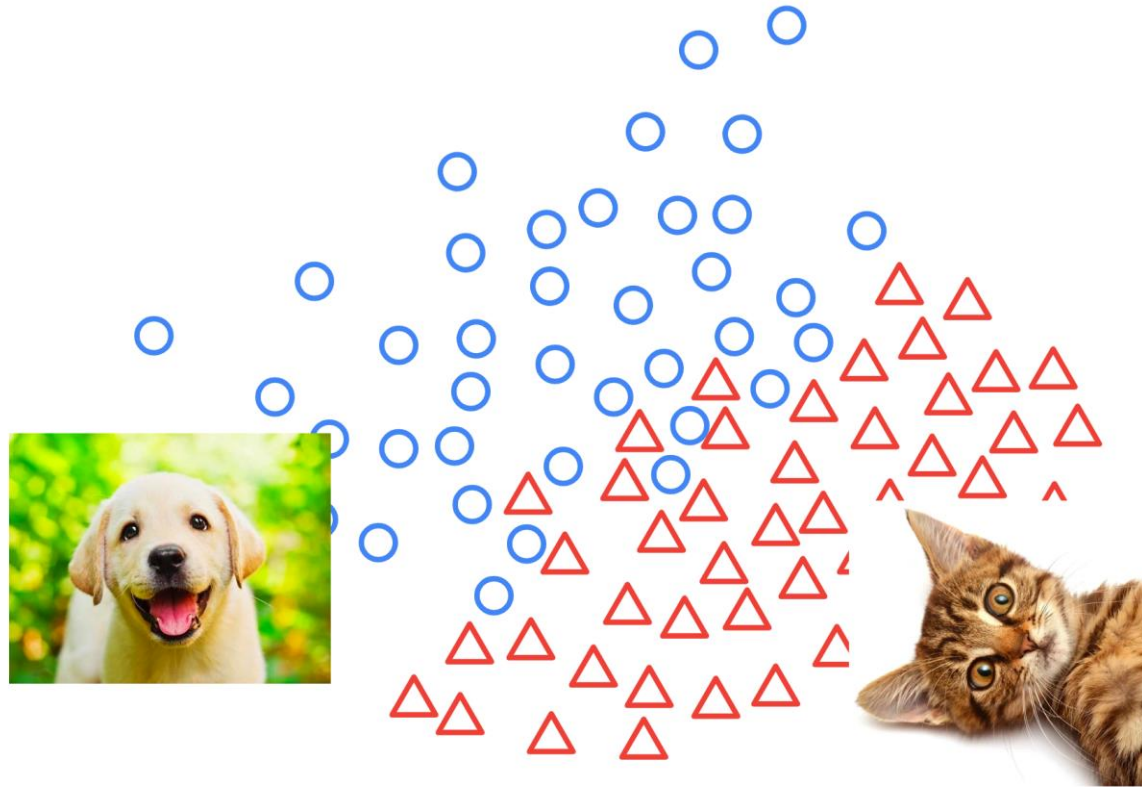
Classification: predict a label with *the confidence*



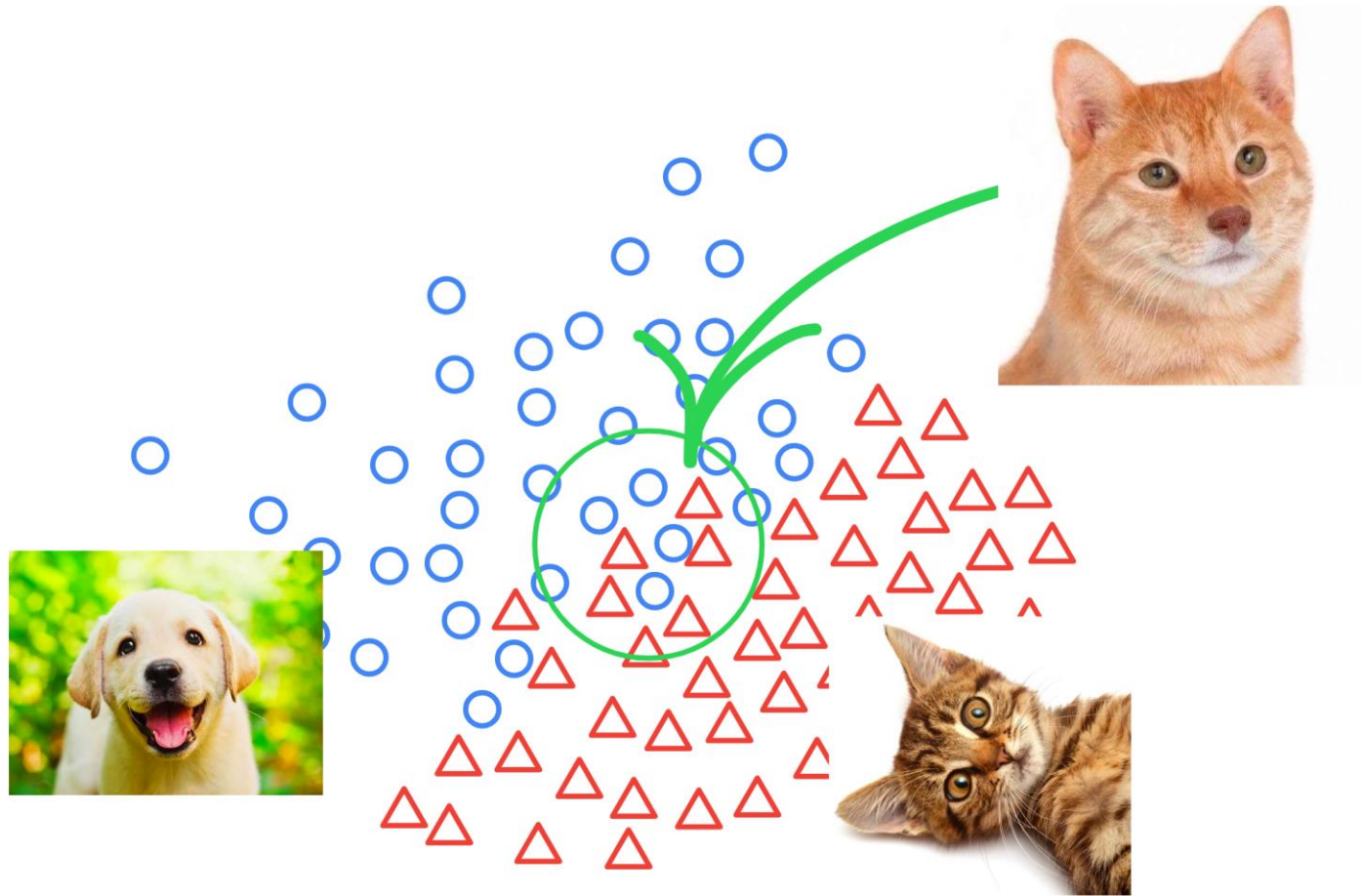
Model/data uncertainties for a regression problem



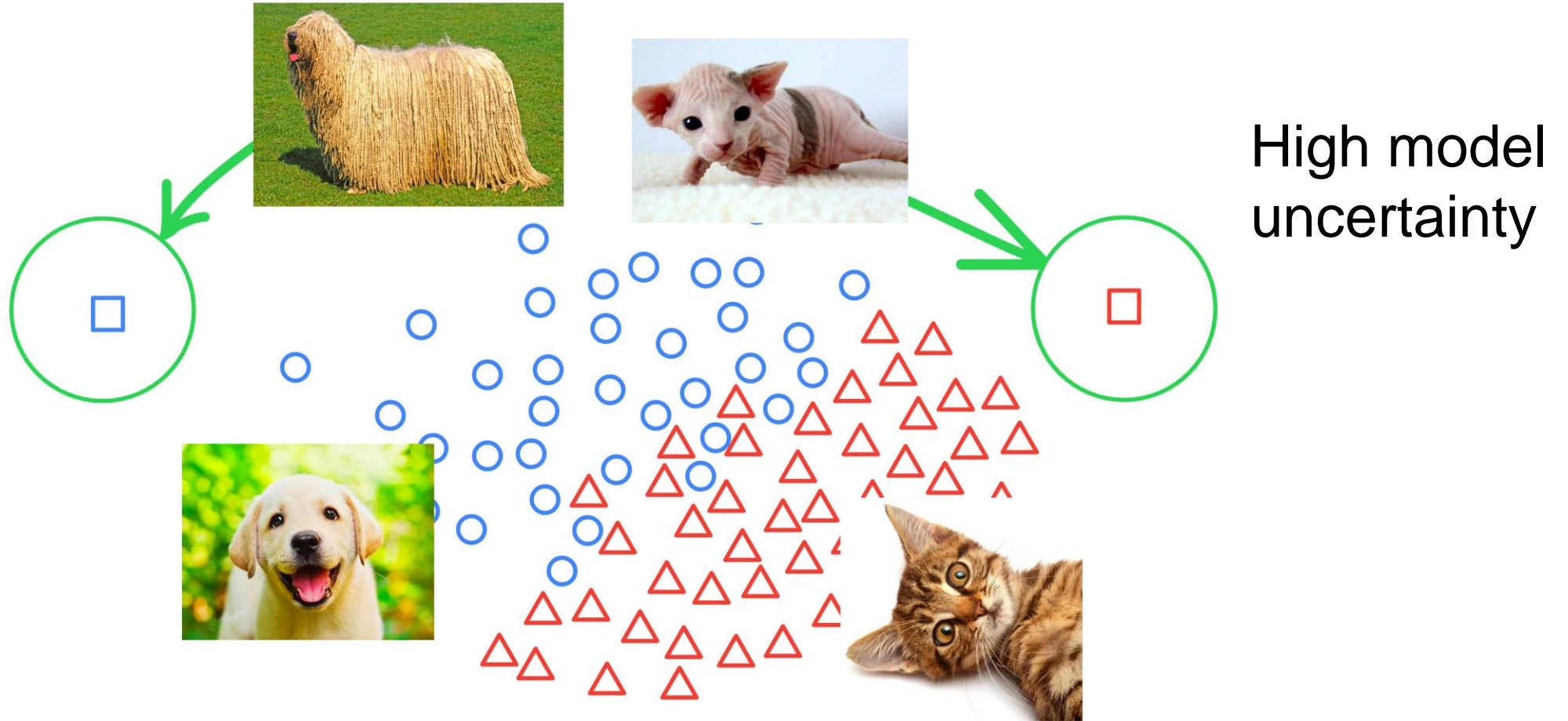
Model/data uncertainties for a classification problem



Model/data uncertainties for a classification problem



Model/data uncertainties for a classification problem



What type of uncertainty do we need to estimate for active learning?

Uncertainty evaluation for linear models

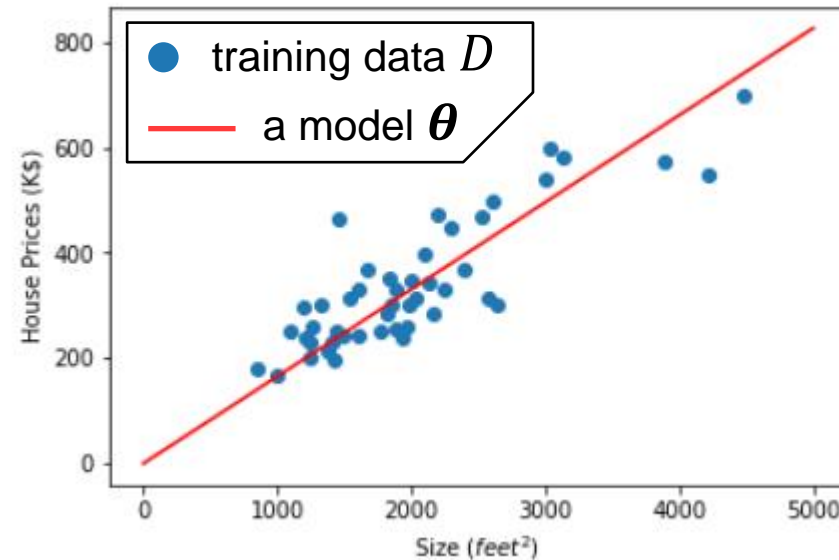
Bayesian linear regression model

We have training data $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n = (X, \mathbf{y})$, $\mathbf{x}_i \in \mathbb{R}^d, y_i \in \mathbb{R}$

Our assumptions about the model:

$$y = f_{\boldsymbol{\theta}}(\mathbf{x}) + \varepsilon = \mathbf{x}^T \boldsymbol{\theta} + \varepsilon, \varepsilon \sim N(0, \sigma^2), \boldsymbol{\theta} \sim N(\mathbf{0}, \sigma_{\boldsymbol{\theta}}^2 \mathbf{I})$$

Example: a house price prediction y via the size of the house in squared feet x



To define the model, we estimate $\boldsymbol{\theta}$ or a distribution of $\boldsymbol{\theta}$ given data D

A posterior distribution of parameters

For linear model $y = \mathbf{x}^T \boldsymbol{\theta} + \varepsilon$ we have prior and posterior distributions

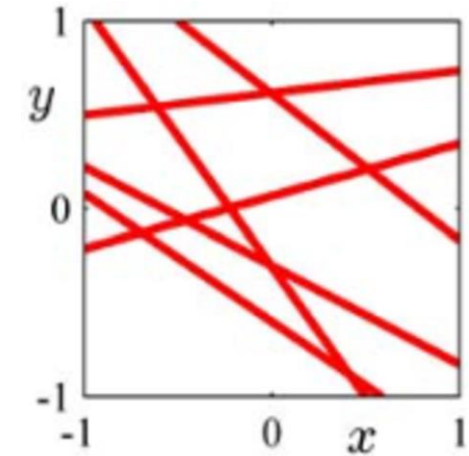
The prior distribution of parameters

$$p(\boldsymbol{\theta}) = N(\mathbf{0}, \sigma_{\theta}^2 \mathbf{I})$$

The posterior distribution for parameters:

$$\begin{aligned} p(\boldsymbol{\theta}|D) &= N(\boldsymbol{\mu}, S), \\ \boldsymbol{\mu} &= \sigma^{-2} S X^T \mathbf{y}, \\ S^{-1} &= \sigma^{-2} X^T X + \sigma_{\theta}^{-2} \mathbf{I} \end{aligned}$$

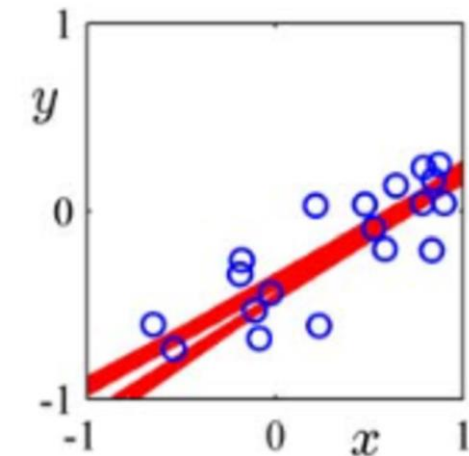
The prior over functions



Training data D



The posterior over functions



Predictive distribution

The predictive distribution corresponds to the predictions at a particular point, we predict not only a mean, but a *distribution*.

The predictive distribution:

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

As both distributions are Gaussian, we can take the integral to get:

$$\begin{aligned} p(y|\mathbf{x}, D) &= N(\mu(\mathbf{x}), \sigma^2(\mathbf{x})), \\ \mu(\mathbf{x}) &= \mathbf{x}^T \boldsymbol{\mu} = \sigma^{-2} \mathbf{x}^T S X^T \mathbf{y}, \\ \sigma^2(\mathbf{x}) &= \mathbf{x}^T (\sigma^{-2} X^T X + \sigma_{\theta}^{-2} \mathbf{I})^{-1} \mathbf{x} + \sigma^2 \end{aligned}$$

Uncertainty estimate for a predictive distribution

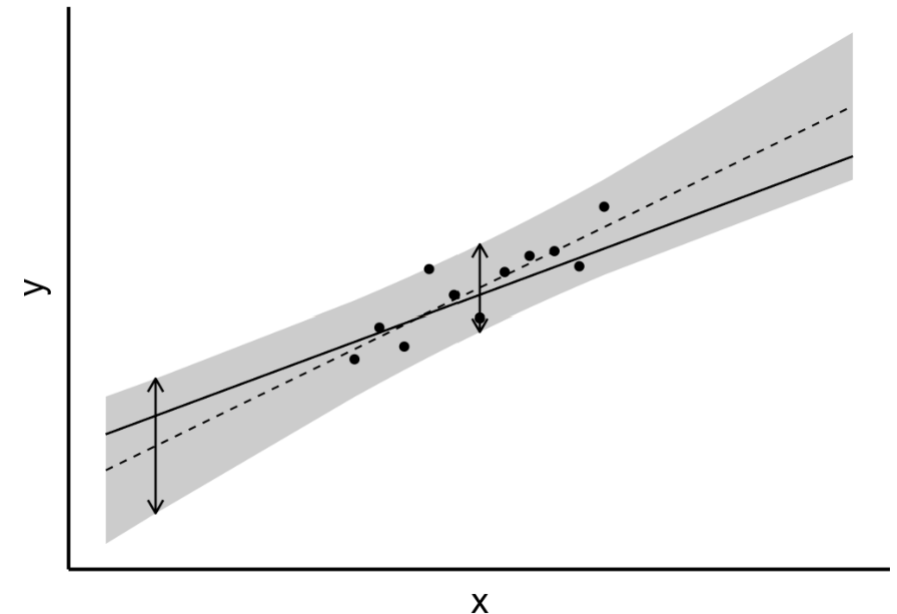
The variance has the analytical form for the linear regression model:

$$\sigma^2(\mathbf{x}) = \underbrace{\mathbf{x}^T (\sigma^{-2} X^T X + \sigma_\theta^{-2} \mathbf{I})^{-1} \mathbf{x}}_{\text{Model uncertainty}} + \underbrace{\sigma^2}_{\text{Data uncertainty}}$$

Skoltech

We solved the problem, but the model is terribly wrong

- Training points
- 95% confidence interval for $N(\mu(\mathbf{x}), \sigma^2(\mathbf{x}))$

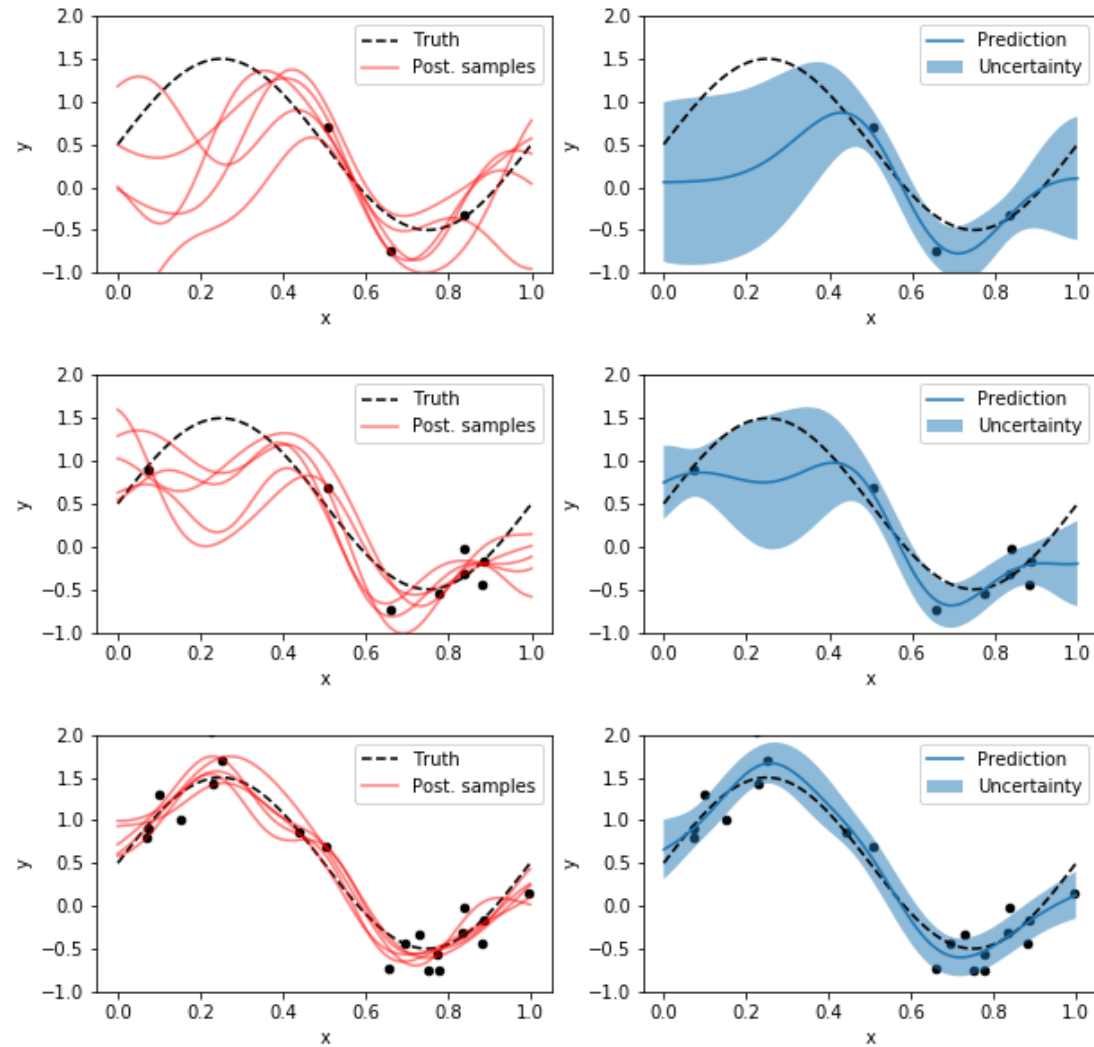


Gaussian process regression

We define a distribution $p(f)$ over a set of functions dense in $C(\mathbb{R}^d)$

In the same way, using the prior distribution $p(f)$ and the data D , we obtain a posterior over functions $p(f|D)$

The predictive distribution has analytical form, if we approximate the posterior $p(f|D)$ with a delta function



Samples from $p(f|D)$

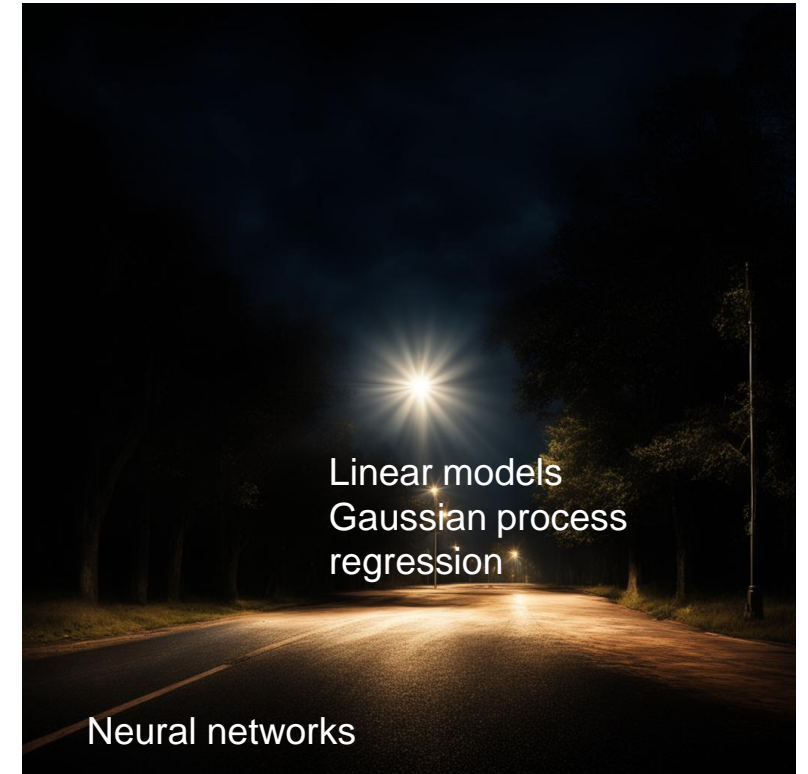
Posterior $p(y|x, \theta, D)$

Good news

- Completely analytical form for the variance
- Rich space of functions

Bad news

- Approximation of the predictive distribution
- We model mostly *model* uncertainties
- **The model is wrong, neural networks work much better**



Streetlight effect for statistics

Pseudo-Bayesian methods via deep ensembles

Why?

In a Bayesian paradigm, we have distributions by construction.

They definitely reflect the uncertainties of predictions

Bayesian Inference for the predictive distribution

In Bayesian paradigm, inference is done via the Posterior Predictive Distribution:

$$p(y \mid \mathbf{x}, D) = \int_{\boldsymbol{\theta}} p(y, \boldsymbol{\theta} \mid \mathbf{x}, D) d\boldsymbol{\theta} = \int_{\boldsymbol{\theta}} p(y \mid \boldsymbol{\theta}, \mathbf{x}) p(\boldsymbol{\theta} \mid D) d\boldsymbol{\theta}$$

However the integral is (almost) impossible to compute for neural networks.

Thus, **approximations should be used.**

Monte-Carlo approximation

To approximate

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

Monte-Carlo estimate is typically used:

$$p(y | \mathbf{x}, D) \approx \frac{1}{K} \sum_{k=1}^K p(y | \boldsymbol{\theta}_k, \mathbf{x}) = \int p(y | \boldsymbol{\theta}, \mathbf{x}) p_K(\boldsymbol{\theta} | D) d\boldsymbol{\theta},$$

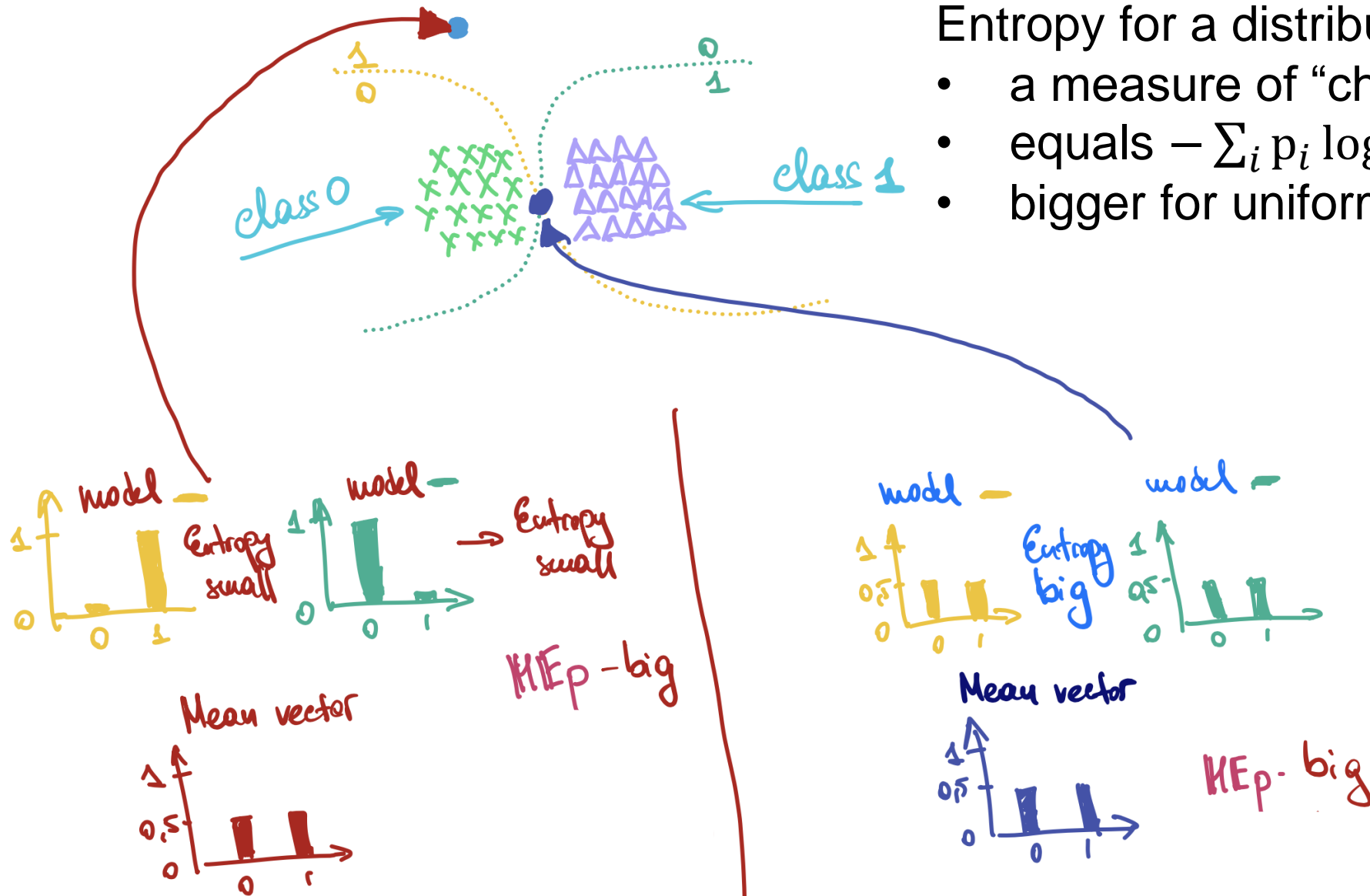
where $\boldsymbol{\theta}_k$ are samples from the posterior distribution over model parameters to from the empirical distribution $p_K(\boldsymbol{\theta} | D)$:

$$\boldsymbol{\theta}_k \sim p(\boldsymbol{\theta} | D)$$

Example for measures of uncertainty

Entropy for a distribution $\{p_i\}_{i=1}^m$:

- a measure of “chaos”
- equals $-\sum_i p_i \log p_i$
- bigger for uniform distributions



Different types of Uncertainty

We derive the following quantities:

Entropy of expected prediction - $\mathbb{H}_p \mathbb{E}_{\theta} p(y \mid \theta, x)$ - **total uncertainty**

Expected entropy - $\mathbb{E}_{\theta} \mathbb{H}_p p(y \mid \theta, x)$ - **data uncertainty**

BALD - $\mathbb{H}_p \mathbb{E}_{\theta} p(y \mid \theta, x) - \mathbb{E}_{\theta} \mathbb{H}_p p(y \mid \theta, x)$ - **model uncertainty**

Note, that these measures can be computed only when we have different samples. Otherwise entropies coincide.

Sampling of models

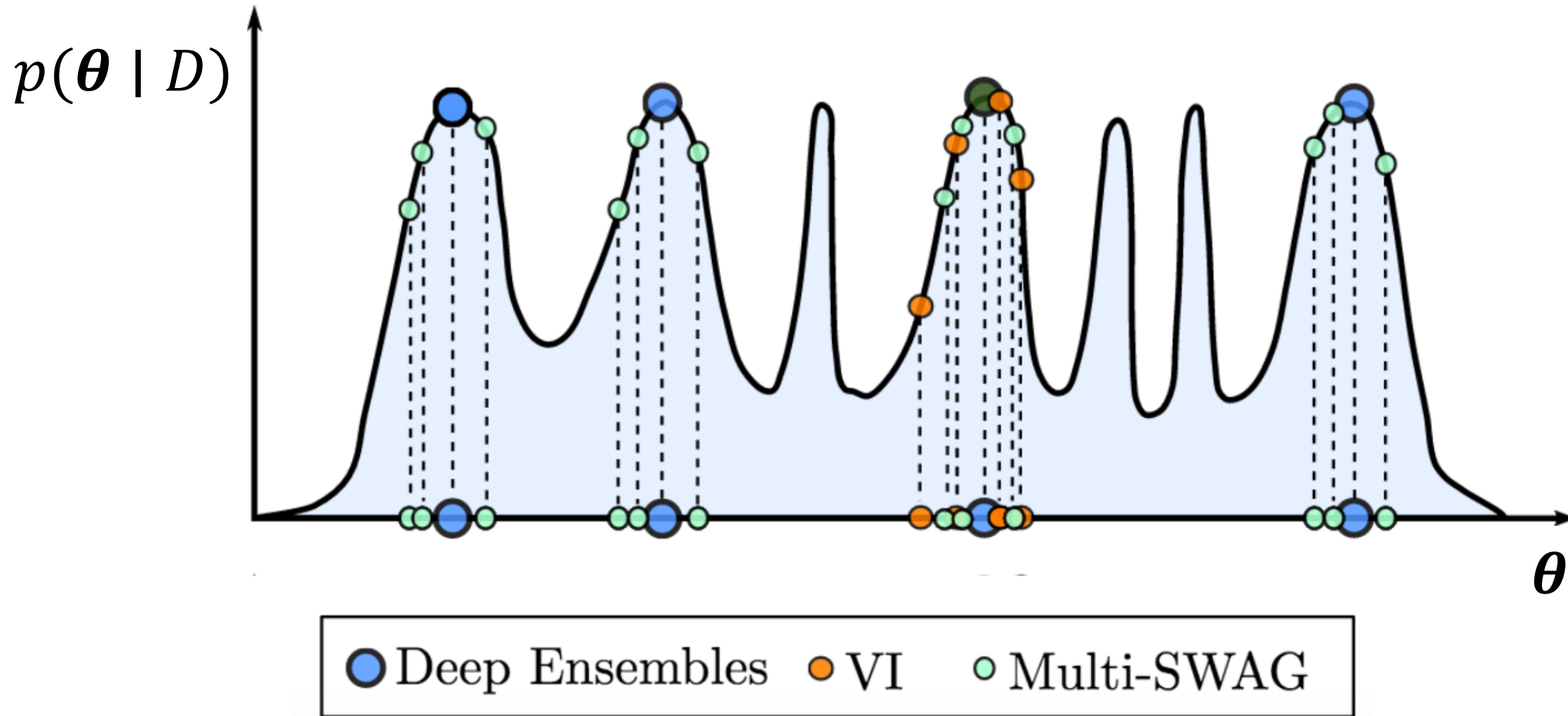
Given samples θ_k , we can approximate all these expectations for uncertainty measures.

But how to receive these samples?

$$\theta_k \sim p(\theta \mid D)$$

- 1. Multistart from different seeds**
- 2. Variational Inference (VI)**
- 3. Markov Chain Monte Carlo (MCMC), SWAG**

Stochastic Weight Averaging Gaussian (SWAG) is a better way to sample an ensemble



Deep Ensembles

Ensemble's pros:

- Accuracy is better than for the single model
- Allows to compute measures of uncertainty

Ensemble's cons:

- Training is K times **more expensive**
- Inference is K times **more expensive**
- Storage is K time **more expensive**

Uncertainty estimation via a single language model: an efficient alternative

Classification baseline, maxprob

It was the best
of times, it was
the worst of...

(1)

Representation



(2)

Classifier

Cats Times

$$\text{uncertainty} = 1 - \max_i p_i$$

Also: entropy
 $-\sum_i p_i \log p_i$

$$\text{certainty} = \max_i p_i$$

$\{p_i\}_{i=1}^m$ m is the dictionary size



Density-based UE Methods

Provide **high-quality UEs**, introduce **low computational overhead**, almost **no additional memory footprint**

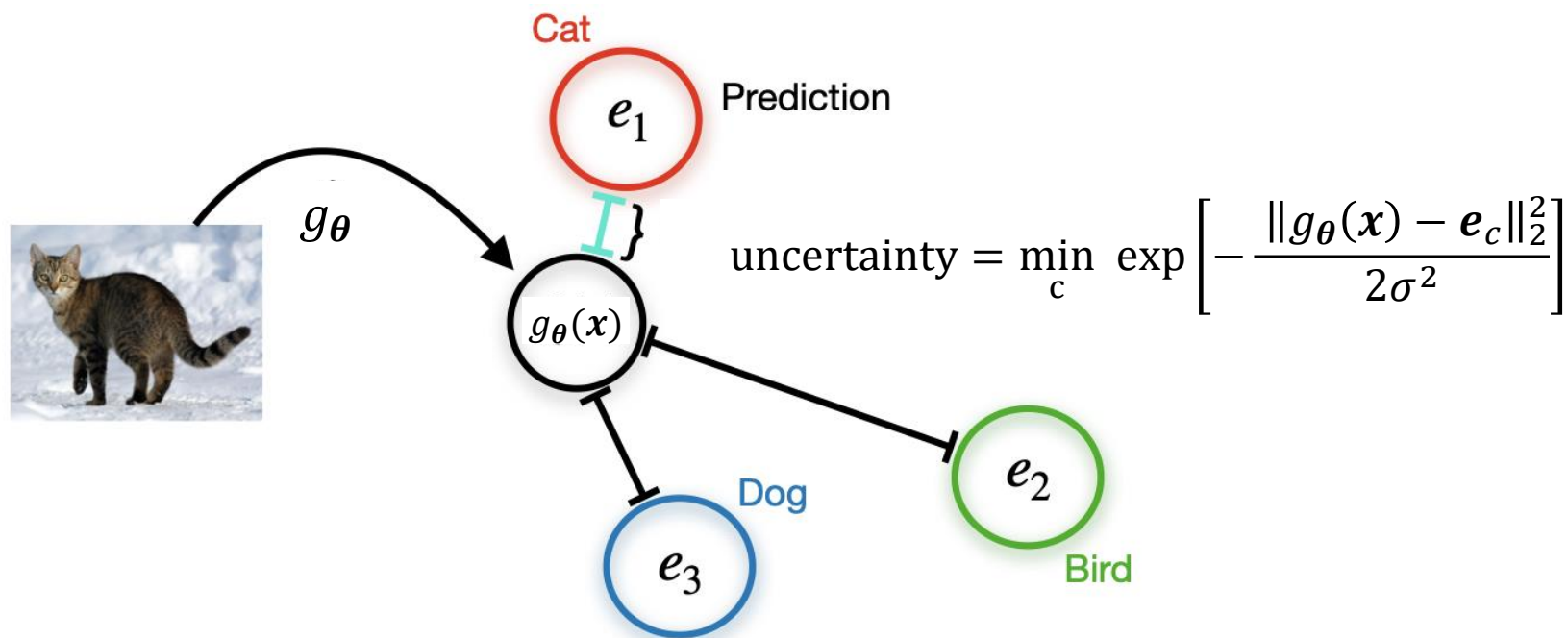
- ↓ Deterministic uncertainty quantification (DUQ) [1]
- ↓ Mahalanobis distance [2]
- ↓ Spectral-normalized Neural Gaussian Process [3]



1. Van Amersfoort, J., et al. "Uncertainty estimation using a single deep deterministic neural network." *ICML*, 2020.
2. Podolskiy, A., et al. "Revisiting mahalanobis distance for transformer-based out-of-domain detection." *AAAI*. 2021.
3. Liu et al. "Simple and Principled Uncertainty Estimation with Deterministic Deep Learning via Distance Awareness." *NeurIPS*. 2020

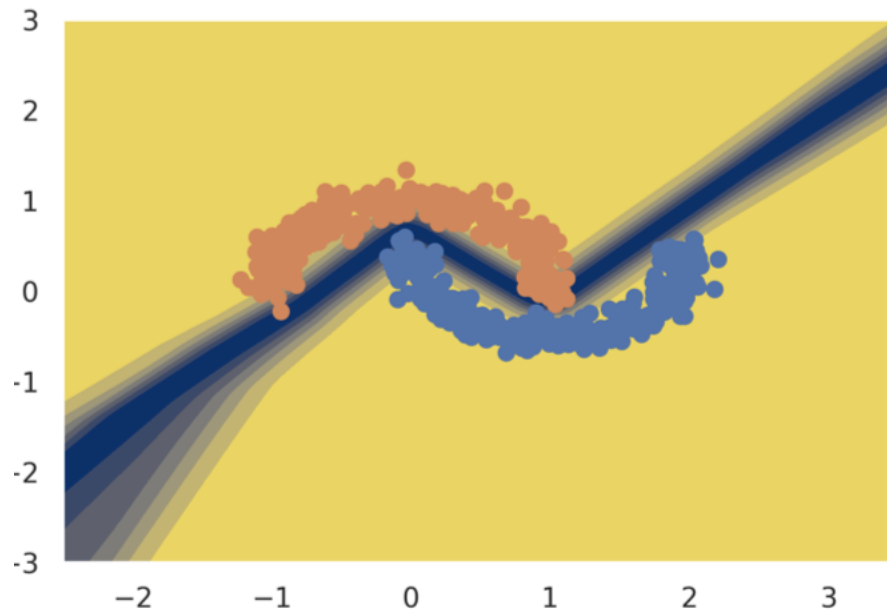
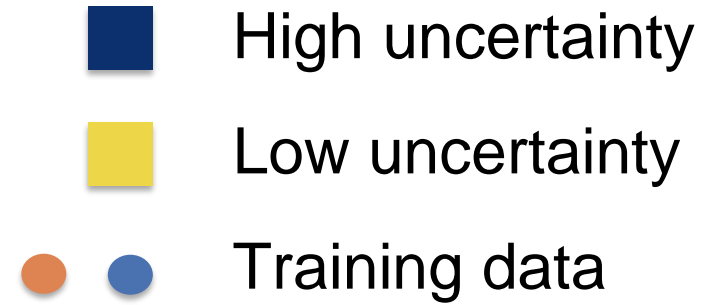
Uncertainty for Deterministic Networks (DUQ)

Common idea — **utilize representations of the training data.**

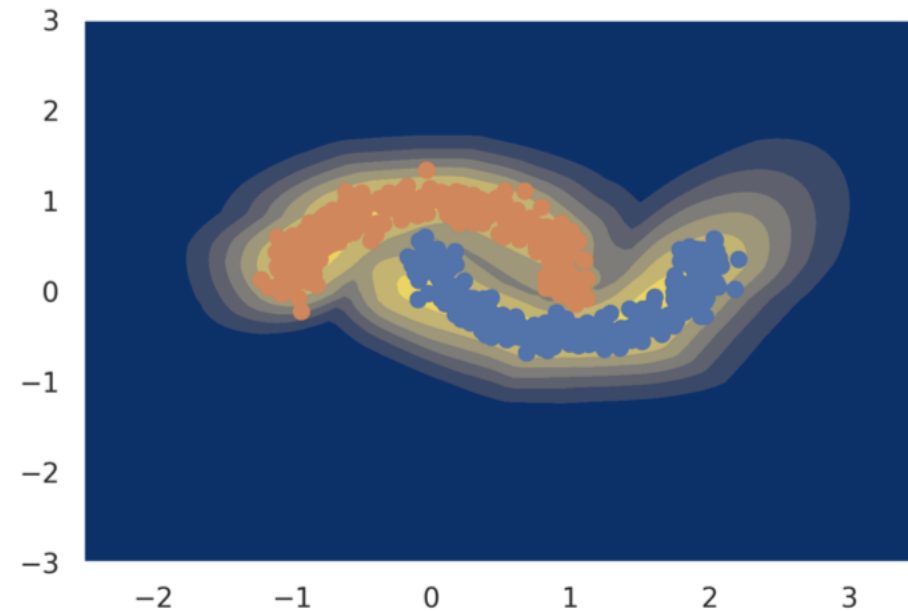


A DUQ architecture. The input is mapped to the feature space, where it is assigned to the closest centroid. The distance to this centroid is uncertainty.

Toy example for DUQ



Deep ensembles



DUQ model

Source: Van Amersfoort, Joost, et al. "Uncertainty estimation using a single deep deterministic neural network." *ICML*, 2020.

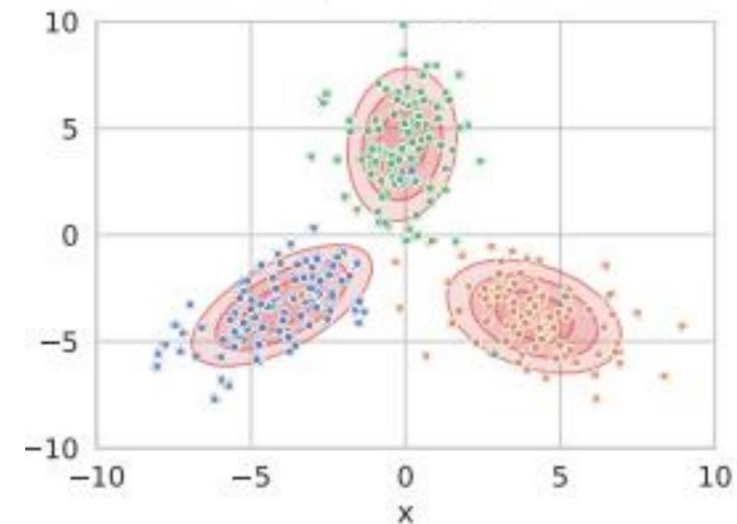
Density-based UE Methods: Deep Deterministic Uncertainty (DDU)

Fit a Gaussian Mixture Model (GMM) on the training data for $p(h(\mathbf{x}))$, where $h(\mathbf{x})$ – hidden representation of instance \mathbf{x} .

$$\bar{U}_E^{DDU}(\mathbf{x}) = \sum_{c \in C} p(h(\mathbf{x})|y = c)p(y = c)$$

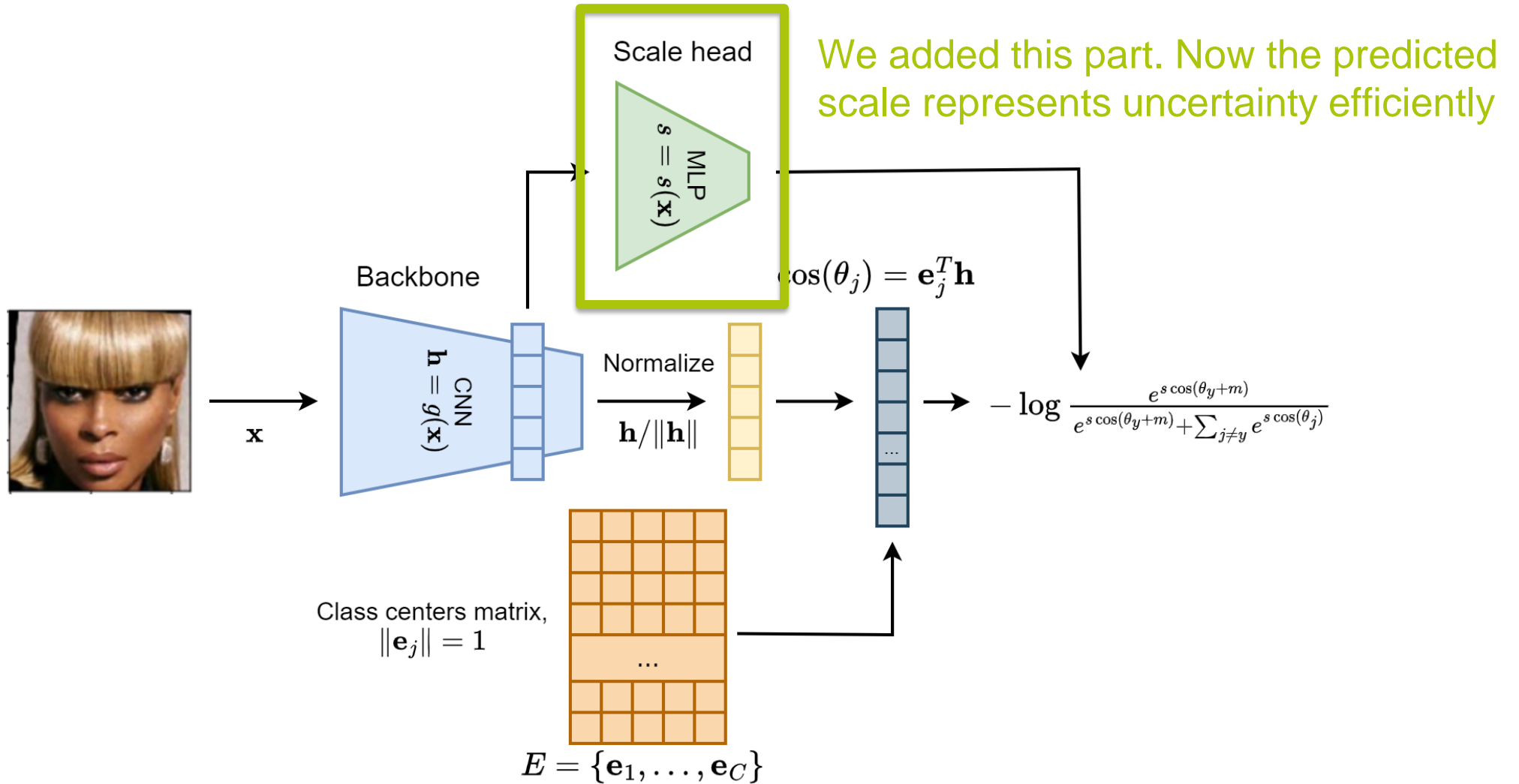
$$p(h(\mathbf{x})|y = c) \sim \mathcal{N}(h(\mathbf{x})|\mu_c, \Sigma_c)$$

$$p(y = c) = \frac{\sum_{(x_i, y_i) \in D} \mathbf{1}[y_i = c]}{|D|}$$

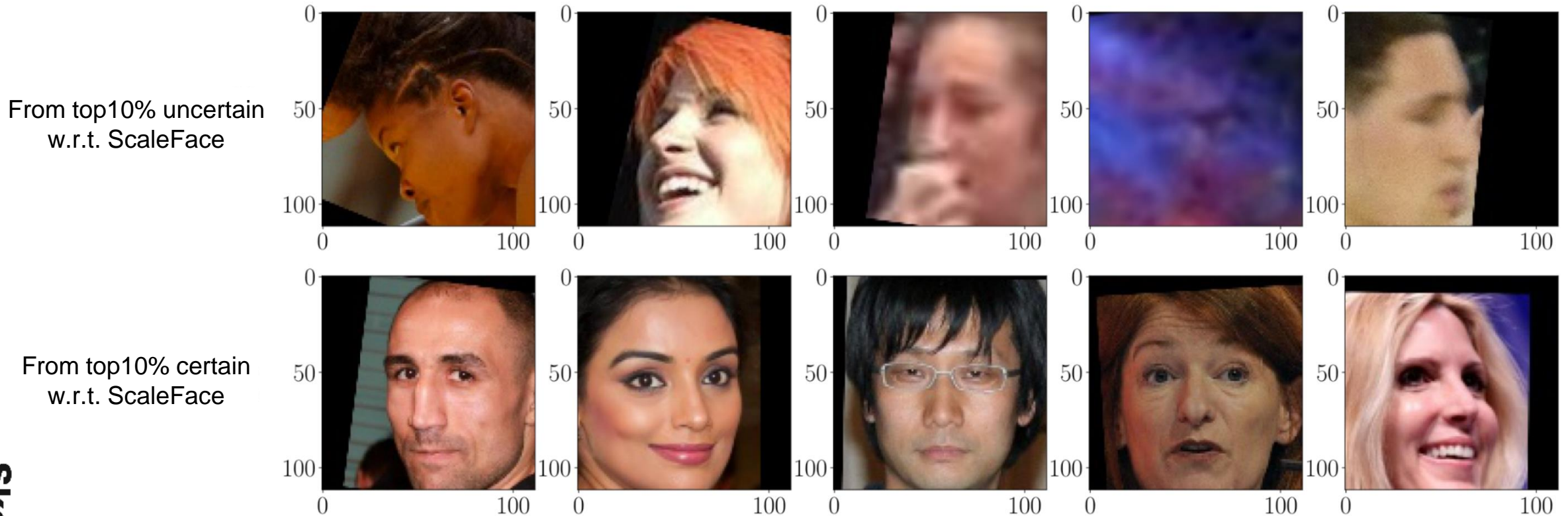


GMM with 3 components fitted to a dataset with 3 different classes

Scale tuning as an alternative



ScaleFace judgements correlate with human for the Face Identification problem

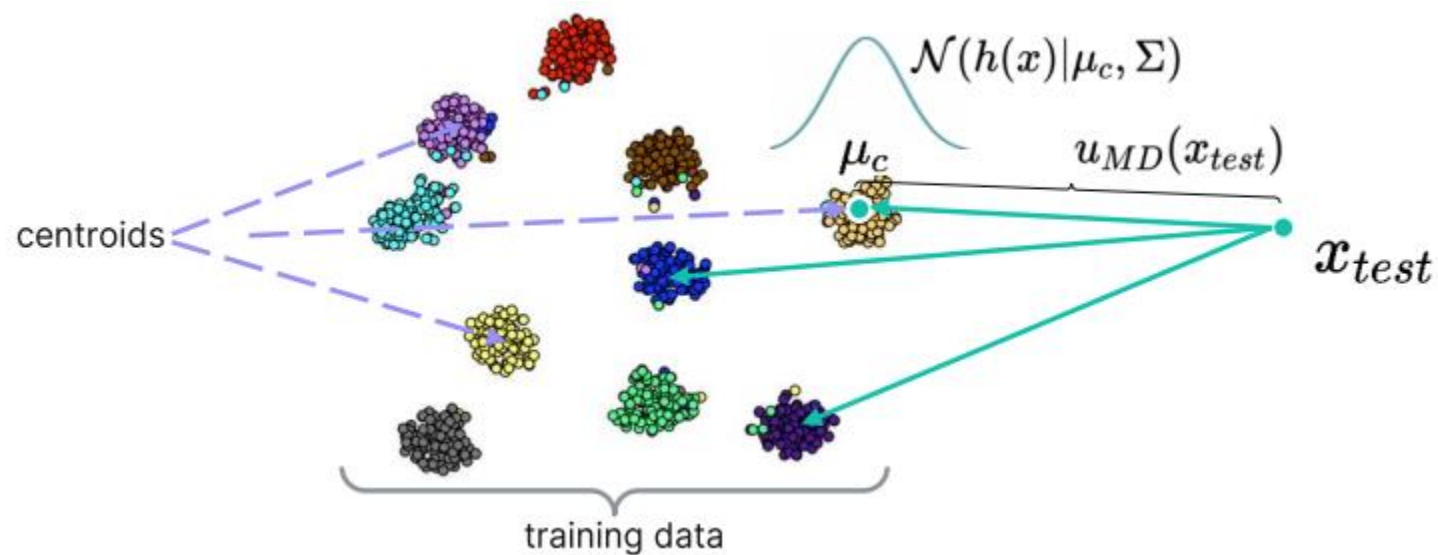


Introducing Mahalanobis Distance

Mahalanobis distance is a generalization of the Euclidian distance. It takes into account the spreading of instances in the training set along various directions in a feature space:

$$u_{MD}(\mathbf{x}_i) = \min_{c \in \mathcal{C}} (\mathbf{h}_i - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}^{-1} (\mathbf{h}_i - \boldsymbol{\mu}_c),$$

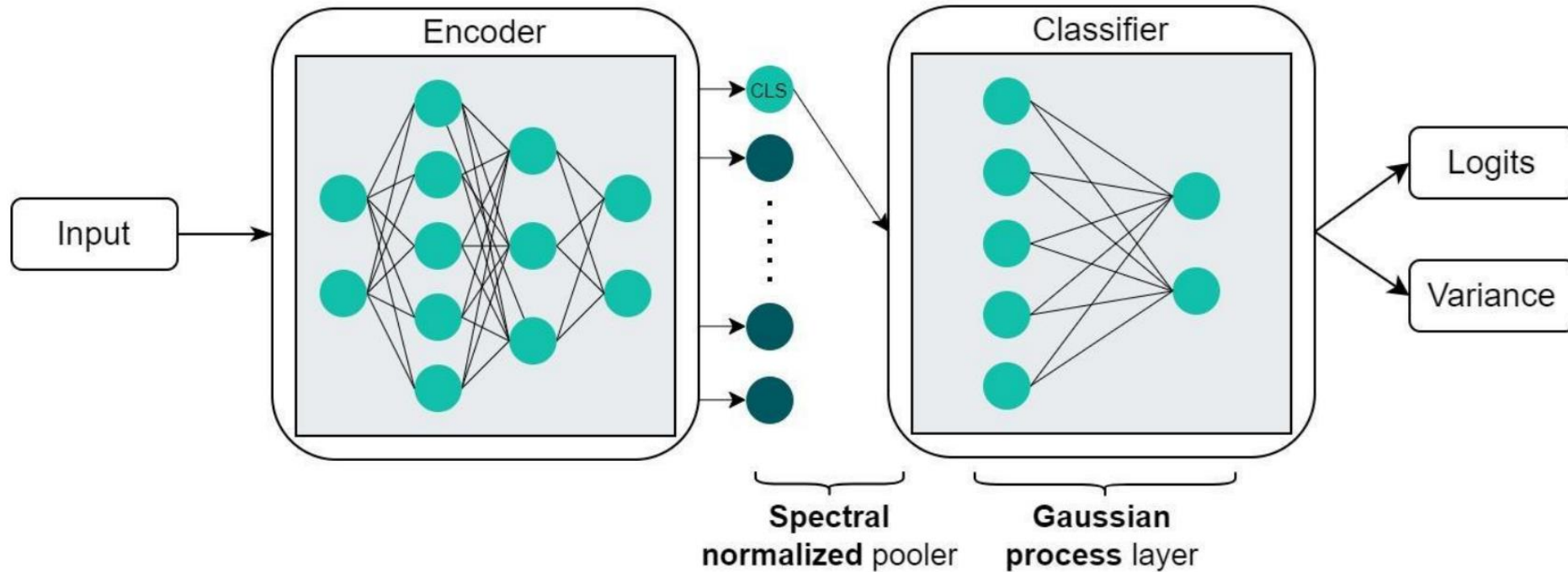
where \mathbf{h}_i is a hidden representation of a i -th instance, $\boldsymbol{\mu}_c$ is a centroid of a class c , and $\boldsymbol{\Sigma}$ is a covariance matrix for hidden representations of training instances.



Spectral-normalized Neural Gaussian Process

A. Use spectral normalization (SN) in the weight matrix of the one-but-last classification layer to preserve distance for hidden representations.

B. Replace the classic dense output layer of a network with a layer that implements a Gaussian process (GP) with an RBF kernel.



A: Spectral normalization

$$L_1 \|\mathbf{x} - \mathbf{x}'\|_X \leq \|h(\mathbf{x}) - h(\mathbf{x}')\|_H \leq L_2 \|\mathbf{x} - \mathbf{x}'\|_X$$

$$W_l = \begin{cases} c \frac{W_l}{\hat{\lambda}}, & \text{if } c < \hat{\lambda} \\ W_l, & \text{otherwise} \end{cases}$$

We want bi-Lipschitz property for embeddings

To ensure this, we limit the norm of weights matrix, $\hat{\lambda} \approx |W_l|_2$

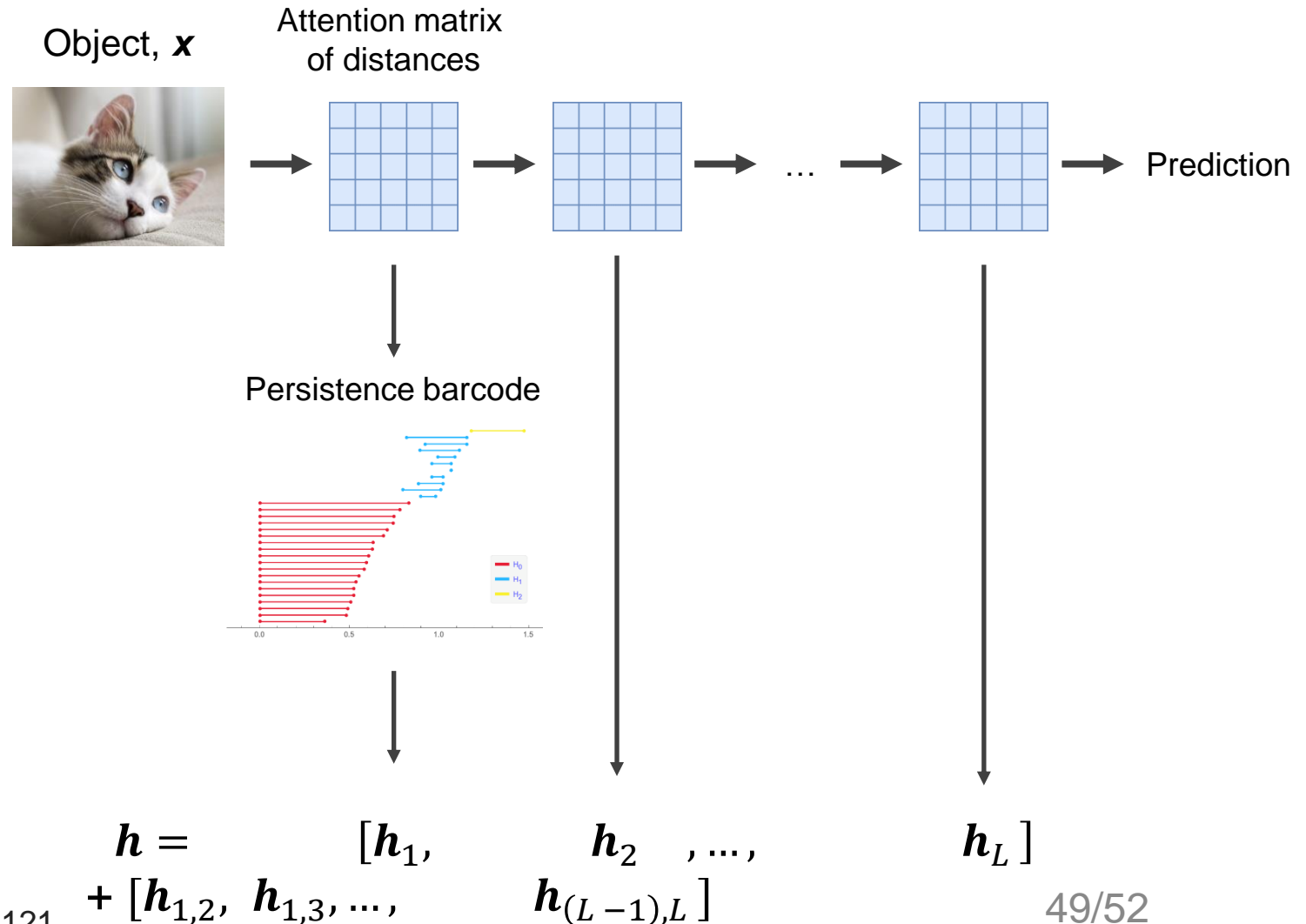
B: Random Fourier Features GP

- Squared exponential kernel
- Random Fourier Features (linear regression in a different feature space)
- Laplace approximation

Topological data analysis (TDA) representations

Existing models are either computationally demanding or look only at the last layer

- TDA helps to extract information about intermediate layers by examining barcodes of attention matrices $[h_1, h_2, \dots, h_L]$ and cross barcodes h_{ij} via **persistent homologies**
- Now we have a representation h that takes into account all layers



Usage of topological data analysis (TDA) for uncertainty estimation

Using the TDA-based representation \mathbf{h} we recover uncertainty better than other methods (1, 2, 3)

We need both common TDA features and features from cross barcodes

Method type	Method name	En-CoLA	Ita-CoLA
Basic Methods	Softmax Response	0.068	0.085
	MC Dropout ¹	0.071	0.084
	Mahalanobis estimator ²	0.083	0.091
	Embedding estimator ³	0.075	0.090
Our methods	Topological estimator without cross-barcodes	<u>0.087</u>	<u>0.092</u>
	Topological estimator with cross-barcodes	0.098	0.099
Oracle Upper Bound		0.124	0.121

Area under rejection curve metric. We want to maximize it. The best value for each dataset (En-CoLA, Ita-CoLA) is highlighted with **bold font**, the second best values – with an underscore.

Semantic uncertainty

- Language models exhibit uncertainty at the token level, where different sentences can convey the same meaning.
- Our primal focus lies on understanding the intended meaning rather than the wording used to convey that meaning.
- Solution – semantic uncertainty

1. **Generation:** Sample M sequences from the predictive distribution of a large language model given a context x .
2. **Clustering:** Cluster the sequences which mean the same thing using bidirectional entailment algorithm.
3. **Entropy estimation:** Approximate semantic entropy by summing probabilities that share a meaning and compute resulting entropy.

(a) Scenario 1: No semantic equivalence

Answer s	Likelihood $p(s x)$	Semantic likelihood $\sum_{s \in c} p(s x)$
Paris	0.5	0.5
Rome	0.4	0.4
London	0.1	0.1
Entropy	0.94	0.94



(b) Scenario 2: Some semantic equivalence

Answer s	Likelihood $p(s x)$	Semantic likelihood $\sum_{s \in c} p(s x)$
Paris	0.5	0.9
It's Paris	0.4	
London	0.1	0.1
Entropy	0.94	0.33

Open problems in uncertainty estimation

1. This area lacks theoretical explanations of observed phenomena.
2. Compared to other methods such as Gaussian process regression, uncertainties for neural networks are inferior.
3. The most effective techniques necessitate the use of ensembles, which leads to computational inefficiency.
4. Furthermore, these methods demonstrate suboptimal performance when applied to Large Language Models.



A musician or a face?