

Active Learning & Beyond

Stefano Teso

University of Trento, Italy

Preliminaries

Strategies

Extensions

Conclusion and Further Reading

Preliminaries

"Imagine that you are the leader of a colonial expedition from Earth to an extrasolar planet. Luckily, this planet is habitable and has a fair amount of vegetation suitable for feeding your group. Importantly, the most abundant source of food comes from a plant whose fruits are sometimes smooth and round, but sometimes bumpy and irregular."

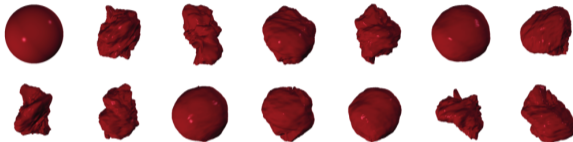


Figure 1.1: Several alien fruits, which vary in shape from round to irregular.

"The physicians assure you that the shape of a fruit is the only feature that seems related to its safety. The problem, though, is that a wide variety of fruit shapes from these plants exist: almost a continuous range from round to irregular. Since the colony has essential uses for both safe and noxious fruits, you want to be able to classify them as accurately as possible. "

Source: [Settles, 2012].

- We know that *smoother* fruits are (monotonically) *safer*, but we don't know where to set the **threshold**.

■ We know that *smoother* fruits are (monotonically) *safer*, but we don't know where to set the **threshold**.

■ In other words, we want to learn a **threshold function**:

$$f_{\theta}(\mathbf{x}) = \begin{cases} 1 & \text{if } x_{\text{irreg}} < \theta \\ -1 & \text{otherwise} \end{cases}$$

where \mathbf{x} are measurements of fruit features and x_{irreg} captures its shape “irregularity”.

Idea: use regular supervised learning

- Collect a large enough training set $S = \{(x, y)\}$, fit threshold classifier f_θ on S
- If maximum % errors is $\epsilon \in (0, 1)$, enough to collect $O(\frac{1}{\epsilon})$ examples [Shalev-Shwartz and Ben-David, 2014]

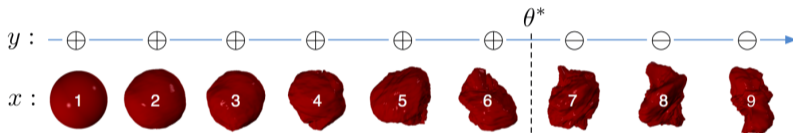


Figure 1.2: Supervised learning for the alien fruits example. Given a set of $\langle x, y \rangle$ instance-label pairs, we want to choose the threshold θ^* that classifies them most accurately.

Idea: use regular supervised learning

- Collect a large enough training set $S = \{(x, y)\}$, fit threshold classifier f_θ on S
- If maximum % errors is $\epsilon \in (0, 1)$, enough to collect $O(\frac{1}{\epsilon})$ examples [Shalev-Shwartz and Ben-David, 2014]

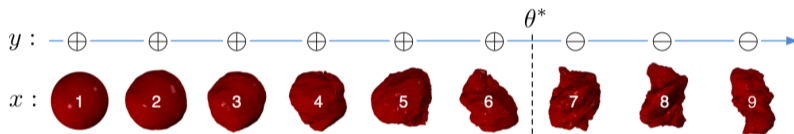


Figure 1.2: Supervised learning for the alien fruits example. Given a set of $\langle x, y \rangle$ instance-label pairs, we want to choose the threshold θ^* that classifies them most accurately.

■ We want to find θ as quickly and as economically as possible, by requiring fewer tests.

■ Can we do better?

Key features:

- Fruits x are plentiful and easy to harvest and measure
- Obtaining y incurs a **cost**: person that eats the fruit may get sick

So we definitely want to **minimize the number of needed labels**.

Idea: gather large set of **unlabeled** fruits $U = \{x_j\}$ and arrange them by roughness.

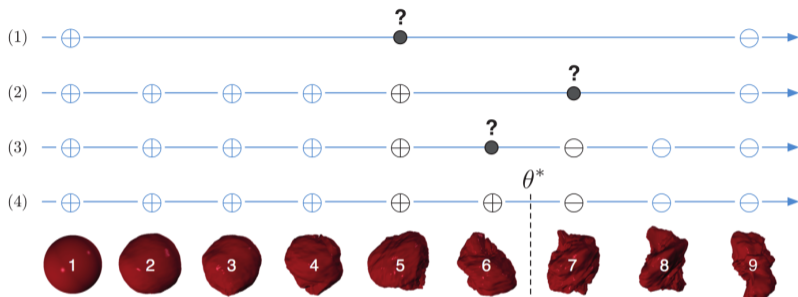


Figure 1.3: A **binary search** through the set of ordered, untested alien fruits. By only testing this subset of fruits, we can exponentially reduce costs while achieving the same result. The labels shown in light blue can be inferred, and therefore do not need to be tested.

Then use **binary search** (as in the illustration) to figure out the threshold θ . This only requires $O(\log_2 \frac{1}{\epsilon})$ tests!

Idea: gather large set of **unlabeled** fruits $U = \{x_i\}$ and arrange them by roughness, then use **binary search**:

ϵ	$\frac{1}{\epsilon}$	$\log_2 \frac{1}{\epsilon}$
0.1	10	3.321
0.001	1000	9.966
0.00001	100000	16.610

■ In this (cleverly designed illustrative) scenario, there is an **exponential** improvement in sample complexity

“The key hypothesis is that if the learner is allowed to choose the data from which it learns — to be active, curious, or exploratory, if you will — it can perform better with less training.” [Settles, 2012]

Preconditions:

- Collecting unlabelled instances x is **cheap**
- Obtaining their labels y is **expensive**

Preconditions:

- Collecting unlabelled instances x is **cheap**
- Obtaining their labels y is **expensive**

Example: Citizen Science

There are tons of images of celestial bodies (think sky surveys). However, in order to understand what's in an image (is it a spiral galaxy? is it a gravitational lensing effect?) you have to ask a human expert.

Preconditions:

- Collecting unlabelled instances x is **cheap**
- Obtaining their labels y is **expensive**

Example: Citizen Science

There are tons of images of celestial bodies (think sky surveys). However, in order to understand what's in an image (is it a spiral galaxy? is it a gravitational lensing effect?) you have to ask a human expert.

Example: Recommendation

There are millions of products on online catalogues (think Amazon), but in order to discover what are the tastes of a user, you have to actually convince them to score the items. This information is personalized, so this is the only way to obtain supervision.

- Adam, the “robot scientist” [King et al., 2009]

The Automation of Science

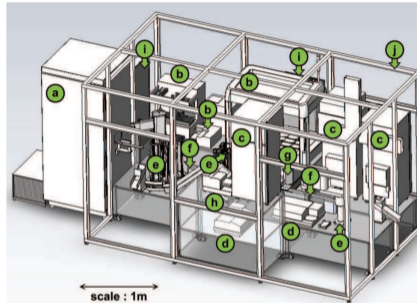
Ross D. King,^{1*} Jem Rowland,¹ Stephen G. Oliver,² Michael Young,³ Wayne Aubrey,¹
Emma Byrne,¹ Maria Liakata,¹ Magdalena Markham,¹ Pinar Pir,² Larisa N. Soldatova,¹
Andrew Sparkes,¹ Kenneth E. Whelan,¹ Amanda Clare¹

The basis of science is the hypothetico-deductive method and the recording of experiments in sufficient detail to enable reproducibility. We report the development of Robot Scientist “Adam,” which advances the automation of both. Adam has autonomously generated functional genomics hypotheses about the yeast *Saccharomyces cerevisiae* and experimentally tested these hypotheses by using laboratory automation. We have confirmed Adam’s conclusions through manual experiments. To describe Adam’s research, we have developed an ontology and logical language. The resulting formalization involves over 10,000 different research units in a nested treelike structure, 10 levels deep, that relates the 6.6 million biomass measurements to their logical description. This formalization describes how a machine contributed to scientific knowledge.

- The learner obtains labels by operating an automated testing machine.

Example: Scientific Discovery

Fig. 1. The Robot Scientist Adam. The advances that distinguish Adam from other complex laboratory systems are the individual design of the experiments to test hypotheses and the utilization of complex internal cycles. Adam's basic operations are selection of specified yeast strains from a library held in a freezer, inoculation of these strains into microtiter plate wells containing rich medium, measurement of growth curves on rich medium, harvesting of a defined quantity of cells from each well, inoculation of these cells into wells containing defined media (minimal synthetic dextrose medium plus up to four added metabolites from a choice of six), and measurement of growth curves on the specified media. To achieve this functionality, Adam has the following components: a, an automated -20°C freezer; b, three liquid handlers (one of which can separately control 96 fluid channels simultaneously); c, three automated $+30^{\circ}\text{C}$ incubators; d, two automated plate readers; e, three robot arms; f, two automated plate slides; g, an automated plate centrifuge; h, an automated plate washer; i, two high-efficiency particulate air filters; and j, a rigid transparent plastic enclosure. There are also two bar code readers, seven cameras, 20 environment sensors, and four personal computers, as well as the software. Adam is capable of designing and initiating over a thousand new



strain and defined-growth-medium experiments each day (from a selection of thousands of yeast strains), with each experiment lasting up to 5 days. The design enables measurement of $\text{OD}_{550\text{nm}}$ for each experiment at least once every 30 min (more often if running at less than full capacity), allowing accurate growth curves to be recorded (typically we take over a hundred measurements a day per well, plus associated metadata. See the supporting online material for pictures and a video of Adam in action.

■ Similar strategies used in chemical engineering, material engineering, etc.

A summary of frequently used terms:

- **Instances** $\mathbf{x} \in \mathbb{R}^d$ are *unlabelled* d -dimensional vectors of observations
- **Examples** $z = (\mathbf{x}, y)$ are instances *annotated* by a label $y \in \{0, 1\}$ or $y \in \{1, \dots, c\}$
- A **classifier** $f : \mathbb{R}^d \rightarrow \{0, 1\}$ maps instances to labels
- $\mathcal{F} = \{f_\theta\}$ is a **family of classifiers** parameterized by θ

■ The meaning of θ depends on the model class, e.g., for neural nets with a fixed architecture, θ represents their weights; for random forests, θ represents the structure and leaves of all trees.

Assumptions

- We assume the data to be distributed according to a **ground-truth data generating process** $p^*(Y, \mathbf{X})$, and use $p^*(Y | \mathbf{X})$ and $p^*(\mathbf{X})$ to indicate the prior and class likelihood of this process, i.e.,

$$p^*(Y, \mathbf{X}) \equiv p^*(Y | \mathbf{X}) \cdot p^*(\mathbf{X}) \quad (1)$$

The data may be corrupted by label noise, but we explicitly avoid adversarial settings.

- We also assume to deal with **probabilistic classifiers**:

$$p_\theta(Y = y | \mathbf{X} = \mathbf{x}) \quad (2)$$

Predictions are computed as:

$$f_\theta(\mathbf{x}) := \operatorname{argmax}_{y=1, \dots, c} p_\theta(Y = y | \mathbf{X} = \mathbf{x}) \quad (3)$$

This includes, e.g., logistic regression, Gaussian processes, and neural nets with a top softmax activation.

- Predictors that only output a **score** (e.g., Support Vector Machines output the distance to the separating hyperplane) can often be adapted (via, e.g., calibration or Platt scaling). This is not necessary to leverage AL.

■ The **labeling oracle** $\text{label} : \mathbb{R}^d \rightarrow \{0, 1\}$ returns:

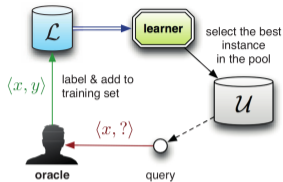
$$\text{label}(\mathbf{x}) := \operatorname{argmax}_{y \in \{0,1\}} p^*(Y = y \mid \mathbf{X} = \mathbf{x}) \quad (4)$$

where $p^*(Y \mid \mathbf{X})$ is the ground-truth distribution.

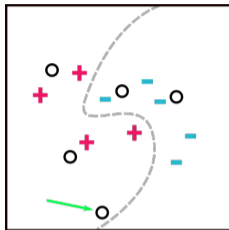
■ Invoking the oracle comes at a **cost**, which is usually non-negligible, instance- and class-dependent, and possibly unknown [Herde et al., 2021].¹

For simplicity, we assume the cost to be **unitary**: querying any instance $\mathbf{x} \in \mathbb{R}^d$ costs the same.

¹Although it could be estimated from interaction data.



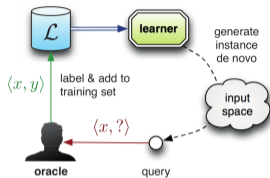
(b) pool-based sampling



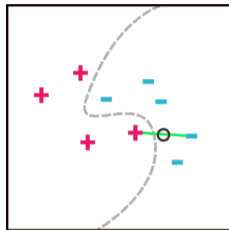
Active Learning (**Pool-based**). **Given:**

- a family of classifiers \mathcal{F} ,
- a set of unlabelled instances $U = \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subseteq \mathbb{R}^d$ sampled i.i.d. from $p^*(\mathbf{X})$,
- a (costly) labeling oracle $\text{label} : \mathbb{R}^d \rightarrow \{0, 1\}$,

Find a classifier $\hat{f} \in \mathcal{F}$ that achieves low risk on $p^*(\mathbf{X}, Y)$ while keeping annot. cost low



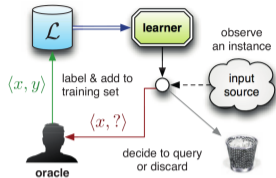
(a) query synthesis



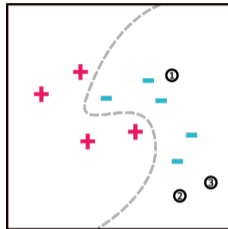
Active Learning (Query Synthesis). **Given:**

- a family of classifiers \mathcal{F} ,
- a generator of instances $\text{synthesize}(\text{region}) \mapsto \mathbf{x}$,
- a (costly) labeling oracle $\text{label} : \mathbb{R}^d \rightarrow \{0, 1\}$,

Find a classifier $\hat{f} \in \mathcal{F}$ that achieves low risk on $p^*(\mathbf{X}, Y)$ while keeping annot. cost low



(a) selective sampling

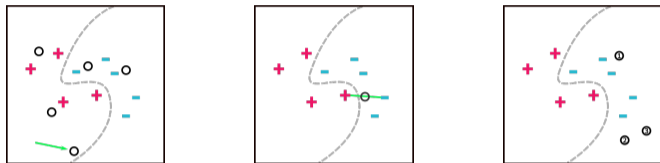


Active Learning (**Selective Sampling**). **Given:**

- a family of classifiers \mathcal{F} ,
- a sequence of unlabelled instances $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots$,
- a (costly) labeling oracle $\text{label} : \mathbb{R}^d \rightarrow \{0, 1\}$

Find a classifier $\hat{f}_t \in \mathcal{F}$ that achieves low risk on future data $\mathbf{x}_{t+1}, \mathbf{x}_{t+2}, \dots$ while keeping annot. cost low

Query Sampling vs. Query Synthesis



■ Left to right:

- **Pool-based:** moderate control over queries, requires memory to store U
- **Query synthesis:** maximum control over queries, can generate uninterpretable queries [Baum and Lang, 1992], although deep generative models can help somehow [Nguyen et al., 2016].
- **Selective sampling:** little control over the distribution of queries, often solved under tight memory constraints (online learning)

■ We will focus on **pool-based AL**.

Strategies

Input: models \mathcal{F} , examples L , pool U , query budget $T \geq 1$

Output: selected model $f \in \mathcal{F}$

```
1:  $f \leftarrow \text{fit}(\mathcal{F}, L)$  ▷ initialize the model
2: for  $t = 1, 2, \dots, T$  do ▷ until the budget is exhausted
3:    $\mathbf{x} \leftarrow \text{argmax}_{\mathbf{x} \in U} \text{acq}(f, \mathbf{x})$  ▷ select a query instance
4:   obtain label  $y$  of  $\mathbf{x}$  from annotator
5:    $U \leftarrow U \setminus \{\mathbf{x}\}$  ▷ remove unlabeled instance from pool
6:    $L \leftarrow L \cup \{(\mathbf{x}, y)\}$  ▷ update training set
7:    $f \leftarrow \text{fit}(\mathcal{F}, L)$  ▷ update the model
return  $f$ 
```

-
- `fit` performs training (e.g., trains for a fixed # of epochs)
 - `acq` scores instances based on their “informativeness”
 - What instance $\mathbf{x} \in U$ should be selected so to convey **as much information as possible** to f ?

What's the point of asking the label of instances on which the classifier is already certain?²

- **Left:** two Gaussians (40 points each)
- **Middle:** picking points completely **at random** (ignoring the class label!)
- **Right:** picking points based on **uncertainty**

²There is a point to doing so, as we will see later.

Uncertainty Sampling

- **Idea:** pick $x \in U$ on which the classifier is **most uncertain**.

Uncertainty Sampling

■ **Idea:** pick $x \in U$ on which the classifier is **most uncertain**.

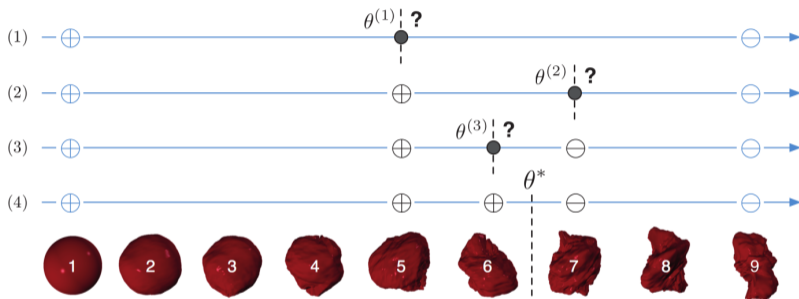


Figure 2.1: The binary search from Figure 1.3, re-interpreted as an uncertainty sampling approach. The best instance to query is deemed to be the one closest to the threshold θ .

Uncertainty Sampling

■ Idea: pick $x \in U$ on which the classifier is **most uncertain**.

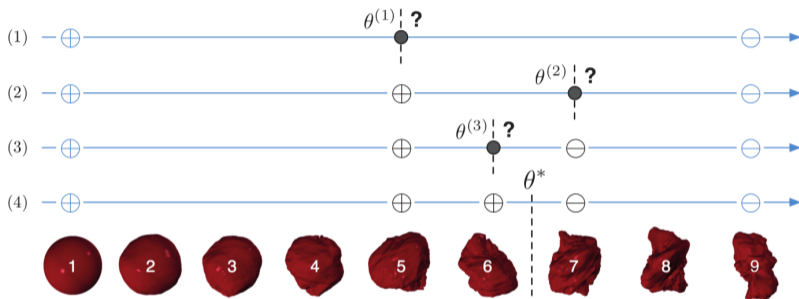


Figure 2.1: The binary search from Figure 1.3, re-interpreted as an uncertainty sampling approach. The best instance to query is deemed to be the one closest to the threshold θ .

■ How should uncertainty be defined?

- Define uncertainty using the **confidence**, i.e., *distance from certainty*:

$$\text{acq}(\theta, \mathbf{x}) := 1 - p_{\theta}(\hat{y} | \mathbf{x}) \quad (5)$$

where \hat{y} is the predicted label:

$$\hat{y} := f_{\theta}(\mathbf{x}) = \underset{y}{\operatorname{argmax}} p_{\theta}(y | \mathbf{x}) \quad (6)$$

- Define uncertainty using the **margin**, i.e., difference in (conditional) likelihood:

$$\text{acq}(\theta, \mathbf{x}) := p_{\theta}(\hat{y}' | \mathbf{x}) - p_{\theta}(\hat{y} | \mathbf{x}) \quad (7)$$

where \hat{y} is the predicted label and \hat{y}' is the **2nd best** label:

$$\hat{y} = \underset{y}{\operatorname{argmax}} p_{\theta}(y | \mathbf{x}) \quad (8)$$

$$\hat{y}' := \underset{y \neq \hat{y}}{\operatorname{argmax}} p_{\theta}(y | \mathbf{x}) \quad (9)$$

- Define uncertainty using the Shannon **entropy** of the label:

$$\text{acq}(\theta, \mathbf{x}) := H_{\theta}(Y | \mathbf{X} = \mathbf{x}) \quad (10)$$

where H_{θ} is defined as:

$$H_{\theta}(Y | \mathbf{X} = \mathbf{x}) := - \sum_{y \in [c]} p_{\theta}(y | \mathbf{x}) \log_2 p_{\theta}(y | \mathbf{x}) \quad (11)$$

- It achieves a **minimum** on dead certain distributions:

$$p_{\theta}(Y | \mathbf{x}) = (0, 1, 0, \dots, 0)$$

- and a **maximum** on the uniform distribution:

$$p_{\theta}(Y | \mathbf{x}) = \left(\frac{1}{c}, \dots, \frac{1}{c}\right)$$

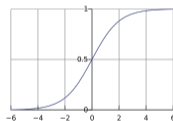
- **Left:** confidence considers prob. of top class only
- **Middle:** margin considers prob. of top & runner up classes
- **Right:** entropy considers prob. of all classes

If $c = 2$, they are equivalent. If $c > 2$, no obvious best choice, it really depends on the task and loss (e.g., cross-entropy vs. accuracy)

- **Left:** confidence considers prob. of top class only
- **Middle:** margin considers prob. of top & runner up classes
- **Right:** entropy considers prob. of all classes

If $c = 2$, they are equivalent. If $c > 2$, no obvious best choice, it really depends on the task and loss (e.g., cross-entropy vs. accuracy)

Example: for classifiers with a sigmoid top layer:



uncertainty depends on distance from separating hyperplane of predicted vs. top two vs. all classes

-
- Uncertainty sampling is very **easy to implement**.
 - Margin & Confidence can be defined even in terms of unnormalized scores.
 - Usually performs reasonably well (though not optimally) in practice: a useful baseline/starting point.

Example: Structured Output

Consider an LSTM that takes a sequence of MNIST images $X = [x_1, \dots, x_n]$ that composes a word and outputs the word itself $y = (y_1, \dots, y_n)$.

- Computing the most likely output \hat{y} can be done efficiently.
- Computing the entropy amounts to:

$$H_{\theta}(Y | X = \mathbf{x}) := - \sum_{y \in \{1, \dots, 26\}^n} p_{\theta}(y | X) \log_2 p_{\theta}(y | X) \quad (12)$$

This involves summing over 26^n possible outputs, which takes time **exponential** in n .

■ Computing the most likely output can be NP-hard. For instance, if y is molecular structure that must satisfy specific hard constraints (chemical validity), then finding the best structure amounts to solving a hard combinatorial problem.

Hence, the **confidence** and **margin** can also be very hard.

- When considering regression models with $Y \in \mathbb{R}$, uncertainty at \mathbf{x} can be implemented as **differential entropy**:

$$H_{\theta}(Y | \mathbf{X} = \mathbf{x}) := \mathbb{E}[-\log_2 p_{\theta}(y | \mathbf{x}) | \mathbf{x}] \quad (13)$$

$$= - \int_{\mathbb{R}} p_{\theta}(y | \mathbf{x}) \log_2 p_{\theta}(y | \mathbf{x}) \quad (14)$$

- When considering regression models with $Y \in \mathbb{R}$, uncertainty at \mathbf{x} can be implemented as **differential entropy**:

$$H_{\theta}(Y | \mathbf{X} = \mathbf{x}) := \mathbb{E}[-\log_2 p_{\theta}(y | \mathbf{x}) | \mathbf{x}] \quad (13)$$

$$= - \int_{\mathbb{R}} p_{\theta}(y | \mathbf{x}) \log_2 p_{\theta}(y | \mathbf{x}) \quad (14)$$

- As an alternative heuristic, use the **variance**:

$$\text{Var}_{\theta}(Y | \mathbf{x}) := \mathbb{E}[(Y - \underbrace{\mathbb{E}[Y | \mathbf{x}]}_{\mu_{\theta}(Y|\mathbf{x})})^2 | \mathbf{x}] \quad (15)$$

$$= \int_{\mathbb{R}} (y - \mu_{\theta}(Y | \mathbf{x}))^2 p_{\theta}(y | \mathbf{x}) dy \quad (16)$$

$$\mu_{\theta}(Y | \mathbf{x}) = \int_{\mathbb{R}} y p_{\theta}(y | \mathbf{x}) dy \quad (17)$$

Uncertainty in Regression

- When considering regression models with $Y \in \mathbb{R}$, uncertainty at \mathbf{x} can be implemented as **differential entropy**:

$$H_{\theta}(Y | \mathbf{X} = \mathbf{x}) := \mathbb{E}[-\log_2 p_{\theta}(y | \mathbf{x}) | \mathbf{x}] \quad (13)$$

$$= - \int_{\mathbb{R}} p_{\theta}(y | \mathbf{x}) \log_2 p_{\theta}(y | \mathbf{x}) \quad (14)$$

- As an alternative heuristic, use the **variance**:

$$\text{Var}_{\theta}(Y | \mathbf{x}) := \mathbb{E}[(Y - \underbrace{\mathbb{E}[Y | \mathbf{x}]}_{\mu_{\theta}(Y|\mathbf{x})})^2 | \mathbf{x}] \quad (15)$$

$$= \int_{\mathbb{R}} (y - \mu_{\theta}(Y | \mathbf{x}))^2 p_{\theta}(y | \mathbf{x}) dy \quad (16)$$

$$\mu_{\theta}(Y | \mathbf{x}) = \int_{\mathbb{R}} y p_{\theta}(y | \mathbf{x}) dy \quad (17)$$

- How to compute them?

Uncertainty in Regression

- Differential entropy and variance:

$$H_{\theta}(Y | X = x) = - \int_{\mathbb{R}} p_{\theta}(y | x) \log_2 p_{\theta}(y | x) \quad \text{Var}_{\theta}(Y | x) = \int_{\mathbb{R}} (p_{\theta}(y | x) - \mu_{\theta}(Y | x))^2 dy \quad (18)$$

³See https://en.wikipedia.org/wiki/Normal_distribution.

Uncertainty in Regression

- Differential entropy and variance:

$$H_{\theta}(Y | \mathbf{X} = \mathbf{x}) = - \int_{\mathbb{R}} p_{\theta}(y | \mathbf{x}) \log_2 p_{\theta}(y | \mathbf{x}) \quad \text{Var}_{\theta}(Y | \mathbf{x}) = \int_{\mathbb{R}} (p_{\theta}(y | \mathbf{x}) - \mu_{\theta}(Y | \mathbf{x}))^2 dy \quad (18)$$

- Both expensive to compute for general models, can **approximate** via quadrature or sampling, but **closed-form solutions** exist for some models (e.g., Gaussian Processes and NNs with a Gaussian output)

³See https://en.wikipedia.org/wiki/Normal_distribution.

Uncertainty in Regression

- Differential entropy and variance:

$$H_{\theta}(Y | \mathbf{X} = \mathbf{x}) = - \int_{\mathbb{R}} p_{\theta}(y | \mathbf{x}) \log_2 p_{\theta}(y | \mathbf{x}) \quad \text{Var}_{\theta}(Y | \mathbf{x}) = \int_{\mathbb{R}} (p_{\theta}(y | \mathbf{x}) - \mu_{\theta}(Y | \mathbf{x}))^2 dy \quad (18)$$

- Both expensive to compute for general models, can **approximate** via quadrature or sampling, but **closed-form solutions** exist for some models (e.g., Gaussian Processes and NNs with a Gaussian output)

Example: 1-dimensional Gaussian Output

Consider one-dimensional output $y \in \mathbb{R}$ and a neural net:

$$nn : \mathbf{x} \mapsto (\mu, \sigma), \quad y \sim \mathcal{N}(\mu, \sigma) \quad (19)$$

In this case, it is well known³ that:

$$\text{Var}_{\theta}(Y | \mathbf{x}) = \sigma^2, \quad H_{\theta}(Y | \mathbf{x}) = \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2} \quad (20)$$

Notice that $\text{Var}_{\theta}(Y | \mathbf{x}) \propto \exp H_{\theta}(Y | \mathbf{x})$, so they change monotonically.

³See https://en.wikipedia.org/wiki/Normal_distribution.

Uncertainty in Regression

- Differential entropy and variance:

$$H_{\theta}(Y | X = \mathbf{x}) = - \int_{\mathbb{R}} p_{\theta}(y | \mathbf{x}) \log_2 p_{\theta}(y | \mathbf{x}) \quad \text{Var}_{\theta}(Y | \mathbf{x}) = \int_{\mathbb{R}} (p_{\theta}(y | \mathbf{x}) - \mu_{\theta}(Y | \mathbf{x})) dy \quad (21)$$

- Both expensive to compute for general models, can **approximate** via quadrature or sampling, but **closed-form solutions** exist for some models (e.g., Gaussian Processes and NNs with a Gaussian output)

Example: k -dimensional Gaussian Output

Consider one-dimensional output $y \in \mathbb{R}^k$ and a neural net:

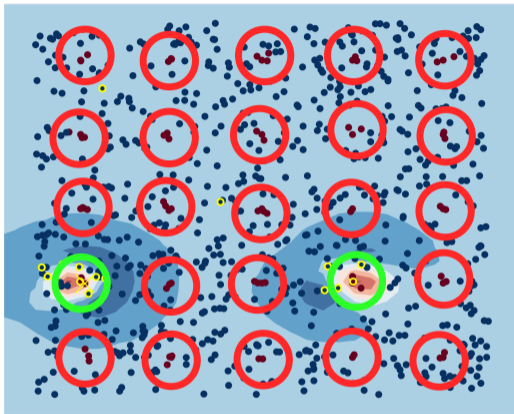
$$nn : \mathbf{x} \mapsto (\boldsymbol{\mu}, S), \quad \boldsymbol{\Sigma} \leftarrow SS^T, \quad \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (22)$$

with $\boldsymbol{\Sigma}$ PSD by construction. In this case, it is well known⁴ that:

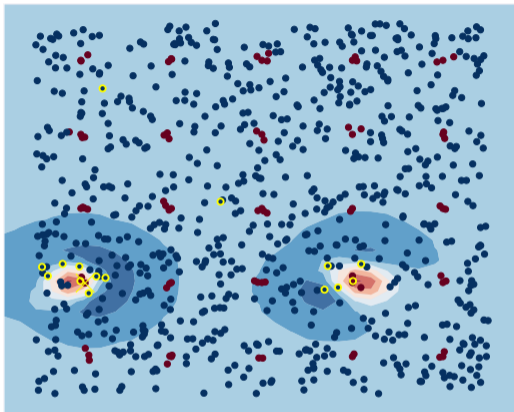
$$\text{Var}_{\theta}(Y | \mathbf{x}) \propto \text{tr} \boldsymbol{\Sigma} \quad H_{\theta}(Y | \mathbf{x}) \propto \log \det \boldsymbol{\Sigma} \quad (23)$$

where the trace is cheap to compute but the determinant is more challenging.

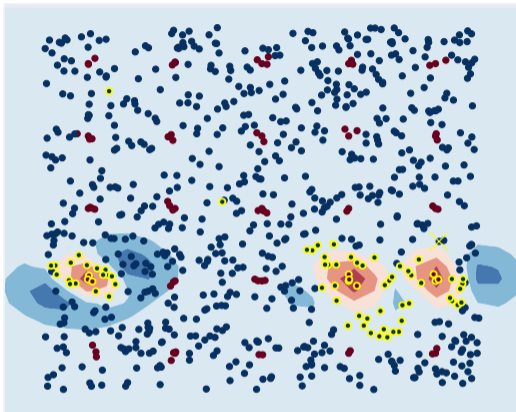
⁴See https://en.wikipedia.org/wiki/Multivariate_normal_distribution.



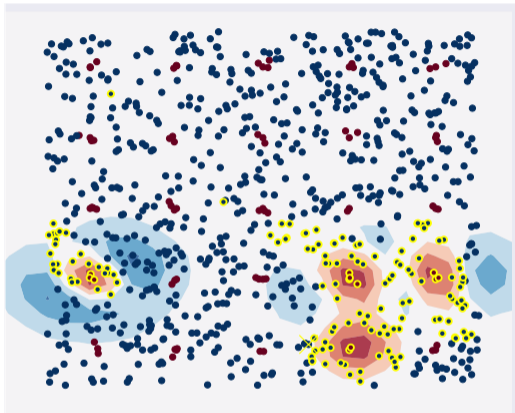
- Synthetic dataset: 25 clusters of **red** points arranged in a 5×5 grid, surrounded by a sea of **blue** points



■ After 10 iterations of uncertainty sampling.



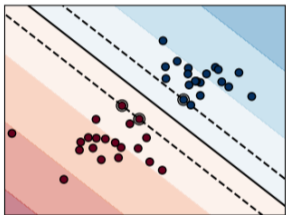
■ After **70** iterations of uncertainty sampling.



■ After **140** iterations of uncertainty sampling. Not nice!

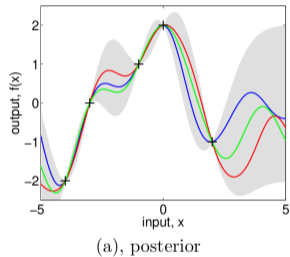
Over-confidence

■ **Discriminative** models are **over-confident**:



Uncertainty does **not** decrease with distance from the training set.

■ **Bayesian generative models** not so much:



Uncertainty **does** decrease with distance from the training set.

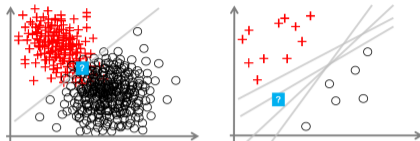


Figure 5: Left: Even with precise knowledge about the optimal hypothesis, the prediction at the query point (indicated by a question mark) is aleatorically uncertain, because the two classes are overlapping in that region. Right: A case of epistemic uncertainty due to a lack of knowledge about the right hypothesis, which is in turn caused by a lack of data.

Aleatoric vs Epistemic [Hüllermeier and Waegeman, 2021]

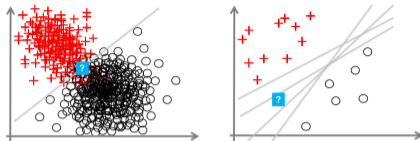


Figure 5: Left: Even with precise knowledge about the optimal hypothesis, the prediction at the query point (indicated by a question mark) is aleatorically uncertain, because the two classes are overlapping in that region. Right: A case of epistemic uncertainty due to a lack of knowledge about the right hypothesis, which is in turn caused by a lack of data.

■ **Aleatoric** uncertainty (“random”) captures how much we can trust the supervision itself. It cannot be decreased. (**left**)

Aleatoric vs Epistemic [Hüllermeier and Waegeman, 2021]

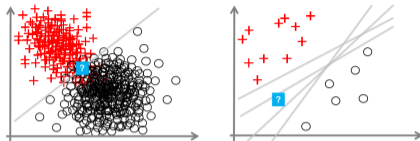


Figure 5: Left: Even with precise knowledge about the optimal hypothesis, the prediction at the query point (indicated by a question mark) is aleatorically uncertain, because the two classes are overlapping in that region. Right: A case of epistemic uncertainty due to a lack of knowledge about the right hypothesis, which is in turn caused by a lack of data.

■ **Aleatoric** uncertainty (“random”) captures how much we can trust the supervision itself. It cannot be decreased. (**left**)

■ **Epistemic** uncertainty (“relating to knowledge”) captures how little we know about the world. This reflects on uncertainty on the choice of θ . It decreases by acquiring **more data**. (**right**)

Aleatoric vs Epistemic [Hüllermeier and Waegeman, 2021]

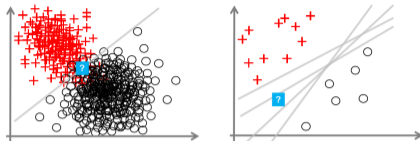


Figure 5: Left: Even with precise knowledge about the optimal hypothesis, the prediction at the query point (indicated by a question mark) is aleatorically uncertain, because the two classes are overlapping in that region. Right: A case of epistemic uncertainty due to a lack of knowledge about the right hypothesis, which is in turn caused by a lack of data.

■ **Aleatoric** uncertainty (“random”) captures how much we can trust the supervision itself. It cannot be decreased. (**left**)

■ **Epistemic** uncertainty (“relating to knowledge”) captures how little we know about the world. This reflects on uncertainty on the choice of θ . It decreases by acquiring **more data**. (**right**)

■ There isn't much point in trying to reduce **aleatoric** uncertainty in AL [Sharma and Bilgic, 2017]

Uncertainty Sampling for Streaming Data

Input: models \mathcal{F} , bootstrap training set L , threshold τ

Output: selected model $f \in \mathcal{F}$

```
1:  $f \leftarrow \text{fit}(\mathcal{F}, L)$  ▷ initialize the model
2: for  $t = 1, 2, 3, \dots$ , do
3:   receive instance  $x$ 
4:   if  $\text{unc}(f, x) > \tau$  then ▷ if  $f$  is uncertain about  $x$ 
5:     obtain label  $y$  of  $x$  from annotator
6:      $L \leftarrow L \cup \{(x, y)\}$  ▷ update training set
7:      $f \leftarrow \text{fit}(\mathcal{F}, L)$  ▷ update the model
return  $f$ 
```

■ The tricky bit is setting τ . Many algorithms update it dynamically by, e.g, starting from a large τ and lowering it as new data is received and the model improves

- For some problems, US **converges** to the right thing – because it is uncertain enough

■ If you are unlucky, US becomes **over-confident**: in this example, the model becomes confident that the regions inside the black blob cannot be white, so it does not sample them and converges to the wrong shape.

- Uncertainty sampling is quite heuristic. Are there more **principled** approaches?

- Consider a **hypothesis space** $\mathcal{F} = \{f_\theta : \mathbb{R}^d \rightarrow [c]\}$ and a data set $S = \{(\mathbf{x}_i, y_i)\}$

- Consider a **hypothesis space** $\mathcal{F} = \{f_\theta : \mathbb{R}^d \rightarrow [c]\}$ and a data set $S = \{(\mathbf{x}_i, y_i)\}$

Consistency

A hypothesis $f \in \mathcal{F}$ is **consistent** with S , written $f \models S$, iff it makes *zero mistakes* on it, that is:

$$(f \models S) \iff \left(\sum_{(\mathbf{x}, y) \in S} \mathbb{1}\{f(\mathbf{x}) \neq y\} \right) = 0 \quad (24)$$

Version Space

- Consider a **hypothesis space** $\mathcal{F} = \{f_\theta : \mathbb{R}^d \rightarrow [c]\}$ and a data set $S = \{(\mathbf{x}_i, y_i)\}$

Consistency

A hypothesis $f \in \mathcal{F}$ is **consistent** with S , written $f \models S$, iff it makes *zero mistakes* on it, that is:

$$(f \models S) \iff \left(\sum_{(\mathbf{x}, y) \in S} \mathbb{1}\{f(\mathbf{x}) \neq y\} \right) = 0 \quad (24)$$

Version Space

The **version space** $VS_{\mathcal{F}}(S)$ of \mathcal{F} given S is the set of hypotheses $f \in \mathcal{F}$ that is *consistent* with S , that is:

$$VS_{\mathcal{F}}(S) = \{f \in \mathcal{F} : f \models S\} \quad (25)$$

Version Space

- Consider a **hypothesis space** $\mathcal{F} = \{f_\theta : \mathbb{R}^d \rightarrow [c]\}$ and a data set $S = \{(\mathbf{x}_i, y_i)\}$

Consistency

A hypothesis $f \in \mathcal{F}$ is **consistent** with S , written $f \models S$, iff it makes *zero mistakes* on it, that is:

$$(f \models S) \iff \left(\sum_{(\mathbf{x}, y) \in S} \mathbb{1}\{f(\mathbf{x}) \neq y\} \right) = 0 \quad (24)$$

Version Space

The **version space** $VS_{\mathcal{F}}(S)$ of \mathcal{F} given S is the set of hypotheses $f \in \mathcal{F}$ that is *consistent* with S , that is:

$$VS_{\mathcal{F}}(S) = \{f \in \mathcal{F} : f \models S\} \quad (25)$$

- $VS_{\mathcal{F}}(L)$ contains those hypotheses that have not yet been ruled out by the acquired examples L : **they are all just as good w.r.t. S**

Version Space

The **version space** $VS_{\mathcal{F}}(S)$ of \mathcal{F} given S is the set of hypotheses $f \in \mathcal{F}$ that is *consistent* with S , that is:

$$VS_{\mathcal{F}}(S) = \{f \in \mathcal{F} : f \models S\} \quad (26)$$

■ If S is not **separable** w.r.t. \mathcal{F} , i.e., if there is no hypothesis $f \in \mathcal{F}$ that is consistent with it, then the version space is empty ($VS_{\mathcal{F}}(S) = \emptyset$). This happens because:

Version Space

The **version space** $VS_{\mathcal{F}}(S)$ of \mathcal{F} given S is the set of hypotheses $f \in \mathcal{F}$ that is *consistent* with S , that is:

$$VS_{\mathcal{F}}(S) = \{f \in \mathcal{F} : f \models S\} \quad (26)$$

■ If S is not **separable** w.r.t. \mathcal{F} , i.e., if there is no hypothesis $f \in \mathcal{F}$ that is consistent with it, then the version space is empty ($VS_{\mathcal{F}}(S) = \emptyset$). This happens because:

- \mathcal{F} is **not expressive enough**. Example: *no $f^* \in \mathcal{F}$ that correctly labels the data.*

Version Space

The **version space** $VS_{\mathcal{F}}(S)$ of \mathcal{F} given S is the set of hypotheses $f \in \mathcal{F}$ that is *consistent* with S , that is:

$$VS_{\mathcal{F}}(S) = \{f \in \mathcal{F} : f \models S\} \quad (26)$$

■ If S is not **separable** w.r.t. \mathcal{F} , i.e., if there is no hypothesis $f \in \mathcal{F}$ that is consistent with it, then the version space is empty ($VS_{\mathcal{F}}(S) = \emptyset$). This happens because:

- \mathcal{F} is **not expressive enough**. Example: *no $f^* \in \mathcal{F}$ that correctly labels the data.*
- S is **noisy**. Example: *S contains the same instance twice but annotated with different labels – e.g., $(x, 1)$ and $(x, 3)$ – so no $f \in \mathcal{F}$ can classify both correctly.*

Version Space

The **version space** $VS_{\mathcal{F}}(S)$ of \mathcal{F} given S is the set of hypotheses $f \in \mathcal{F}$ that is *consistent* with S , that is:

$$VS_{\mathcal{F}}(S) = \{f \in \mathcal{F} : f \models S\} \quad (26)$$

■ If S is not **separable** w.r.t. \mathcal{F} , i.e., if there is no hypothesis $f \in \mathcal{F}$ that is consistent with it, then the version space is empty ($VS_{\mathcal{F}}(S) = \emptyset$). This happens because:

- \mathcal{F} is **not expressive enough**. Example: *no $f^* \in \mathcal{F}$ that correctly labels the data.*
- S is **noisy**. Example: *S contains the same instance twice but annotated with different labels – e.g., $(x, 1)$ and $(x, 3)$ – so no $f \in \mathcal{F}$ can classify both correctly.*

■ We assume the **realizable case**: $\exists f^* \in \mathcal{F}$ s.t. $y = f^*(x)$ for all x and no noise.

This implies that $f^* \in VS_{\mathcal{F}}(L)$ for all choices of labeled examples L , because the supervision (x, y) is always consistent with f^* . In addition, $VS_{\mathcal{F}}(L) \neq \emptyset$.

Disagreement Region

Given \mathcal{F} and S , the **disagreement region** is the set of points $\mathbf{x} \in \mathbb{R}^d$ such that there exist two classifiers f, f' in the version space $VS_{\mathcal{F}}(S)$ that produce different predictions for them:

$$DIS_{\mathcal{F}}(S) = \{\mathbf{x} \in \mathbb{R}^d : \exists f, f' \in VS_{\mathcal{F}}(S) . f(\mathbf{x}) \neq f'(\mathbf{x})\} \quad (27)$$

Disagreement Region

Given \mathcal{F} and S , the **disagreement region** is the set of points $\mathbf{x} \in \mathbb{R}^d$ such that there exist two classifiers f, f' in the version space $VS_{\mathcal{F}}(S)$ that produce different predictions for them:

$$DIS_{\mathcal{F}}(S) = \{\mathbf{x} \in \mathbb{R}^d : \exists f, f' \in VS_{\mathcal{F}}(S) . f(\mathbf{x}) \neq f'(\mathbf{x})\} \quad (27)$$

- If $\mathbf{x} \notin DIS_{\mathcal{F}}(S)$, then all candidate classifiers f in the version space classify it the same: acquiring its label is **pointless**.

Disagreement Region

Given \mathcal{F} and S , the **disagreement region** is the set of points $\mathbf{x} \in \mathbb{R}^d$ such that there exist two classifiers f, f' in the version space $VS_{\mathcal{F}}(S)$ that produce different predictions for them:

$$DIS_{\mathcal{F}}(S) = \{\mathbf{x} \in \mathbb{R}^d : \exists f, f' \in VS_{\mathcal{F}}(S) . f(\mathbf{x}) \neq f'(\mathbf{x})\} \quad (27)$$

- If $\mathbf{x} \notin DIS_{\mathcal{F}}(S)$, then all candidate classifiers f in the version space classify it the same: acquiring its label is **pointless**.
- If $\mathbf{x} \in DIS_{\mathcal{F}}(S)$, then at least one f in the version space classifies it differently: acquiring its label is **useful**.

Left: input space \mathbb{R}^d , data set S of red crosses vs blue circles. **Right:** hypothesis space \mathcal{F} , each f is a point; the ground-truth f^* is in red.

Left: **input space** \mathbb{R}^d , data set S of red crosses vs blue circles. **Right:** **hypothesis space** \mathcal{F} , each f is a point; the ground-truth f^* is in red.

\mathcal{F} is the set of 2D rectangles. Rectangles in instance space (left) are points in hypothesis space (right), as shown by the arrows.

Left: **input space** \mathbb{R}^d , data set S of red crosses vs blue circles. **Right:** **hypothesis space** \mathcal{F} , each f is a point; the ground-truth f^* is in red.

\mathcal{F} is the set of 2D rectangles. Rectangles in instance space (left) are points in hypothesis space (right), as shown by the arrows.

The **version space** $VS_{\mathcal{F}}(S)$ contains all the rectangles (pale gray) between inner & outer rectangles (darker gray)

Left: **input space** \mathbb{R}^d , data set S of red crosses vs blue circles. **Right:** **hypothesis space** \mathcal{F} , each f is a point; the ground-truth f^* is in red.

\mathcal{F} is the set of 2D rectangles. Rectangles in instance space (left) are points in hypothesis space (right), as shown by the arrows.

The **version space** $VS_{\mathcal{F}}(S)$ contains all the rectangles (pale gray) between inner & outer rectangles (darker gray)

The **disagreement region** $DIS_{\mathcal{F}}(S)$ is the space enclosed between these two rectangles.

Input: models \mathcal{F}

Output: selected model $f \in \mathcal{F}$

1: $L \leftarrow \emptyset$

2: $\mathcal{V} \leftarrow \mathcal{F}$

▷ implements the version space $VS_{\mathcal{F}}(L)$

3: **for** $t = 1, 2, 3, \dots$, **do**

4: receive instance \mathbf{x}

5: **if** $\mathbf{x} \in DIS(\mathcal{V})$ **then**

▷ if \mathbf{x} falls in the disagreement region

6: obtain label y of \mathbf{x}

7: update $\mathcal{V} \leftarrow \{f \in \mathcal{V} : f(\mathbf{x}) = y\}$

▷ update version space

8: **return** any $f \in \mathcal{V}$

Input: models \mathcal{F}

Output: selected model $f \in \mathcal{F}$

1: $L \leftarrow \emptyset$

2: $\mathcal{V} \leftarrow \mathcal{F}$

▷ implements the version space $VS_{\mathcal{F}}(L)$

3: **for** $t = 1, 2, 3, \dots$, **do**

4: receive instance x

5: **if** $x \in DIS(\mathcal{V})$ **then**

▷ if x falls in the disagreement region

6: obtain label y of x

7: update $\mathcal{V} \leftarrow \{f \in \mathcal{V} : f(x) = y\}$

▷ update version space

8: **return** any $f \in \mathcal{V}$

■ If $x \in DIS_{\mathcal{F}}(L)$, then there are two $f, f' \in VS_{\mathcal{F}}(L)$ that disagree on how x should be labeled. Getting its label allows us to get rid of either f or f' , so $VS_{\mathcal{F}}(S)$ and $DIS_{\mathcal{F}}(S)$ both shrink.

Input: models \mathcal{F}

Output: selected model $f \in \mathcal{F}$

1: $L \leftarrow \emptyset$

2: $\mathcal{V} \leftarrow \mathcal{F}$

▷ implements the version space $VS_{\mathcal{F}}(L)$

3: **for** $t = 1, 2, 3, \dots$, **do**

4: receive instance \mathbf{x}

5: **if** $\mathbf{x} \in DIS(\mathcal{V})$ **then**

▷ if \mathbf{x} falls in the disagreement region

6: obtain label y of \mathbf{x}

7: update $\mathcal{V} \leftarrow \{f \in \mathcal{V} : f(\mathbf{x}) = y\}$

▷ update version space

8: **return** any $f \in \mathcal{V}$

■ If $\mathbf{x} \in DIS_{\mathcal{F}}(L)$, then there are two $f, f' \in VS_{\mathcal{F}}(L)$ that disagree on how \mathbf{x} should be labeled. Getting its label allows us to get rid of either f or f' , so $VS_{\mathcal{F}}(S)$ and $DIS_{\mathcal{F}}(S)$ both shrink.

■ Recall that f^* is always compatible with examples (\mathbf{x}, y) , so it is always in $\mathcal{V} \rightarrow$ algorithm *zooms* into it!

Input: models \mathcal{F}

Output: selected model $f \in \mathcal{F}$

1: $L \leftarrow \emptyset$

2: $\mathcal{V} \leftarrow \mathcal{F}$

▷ implements the version space $VS_{\mathcal{F}}(L)$

3: **for** $t = 1, 2, 3, \dots$, **do**

4: receive instance \mathbf{x}

5: **if** $\mathbf{x} \in DIS(\mathcal{V})$ **then**

▷ if \mathbf{x} falls in the disagreement region

6: obtain label y of \mathbf{x}

7: update $\mathcal{V} \leftarrow \{f \in \mathcal{V} : f(\mathbf{x}) = y\}$

▷ update version space

8: **return** any $f \in \mathcal{V}$

■ If $\mathbf{x} \in DIS_{\mathcal{F}}(L)$, then there are two $f, f' \in VS_{\mathcal{F}}(L)$ that disagree on how \mathbf{x} should be labeled. Getting its label allows us to get rid of either f or f' , so $VS_{\mathcal{F}}(S)$ and $DIS_{\mathcal{F}}(S)$ both shrink.

■ Recall that f^* is always compatible with examples (\mathbf{x}, y) , so it is always in $\mathcal{V} \rightarrow$ algorithm *zooms* into it!

■ This algorithm makes **no useless queries!**

Input: models \mathcal{F}

Output: selected model $f \in \mathcal{F}$

1: $L \leftarrow \emptyset$

2: $\mathcal{V} \leftarrow \mathcal{F}$

▷ implements the version space $VS_{\mathcal{F}}(L)$

3: **for** $t = 1, 2, 3, \dots$, **do**

4: receive instance \mathbf{x}

5: **if** $\mathbf{x} \in DIS(\mathcal{V})$ **then**

▷ if \mathbf{x} falls in the disagreement region

6: obtain label y of \mathbf{x}

7: update $\mathcal{V} \leftarrow \{f \in \mathcal{V} : f(\mathbf{x}) = y\}$

▷ update version space

8: **return** any $f \in \mathcal{V}$

■ If $\mathbf{x} \in DIS_{\mathcal{F}}(L)$, then there are two $f, f' \in VS_{\mathcal{F}}(L)$ that disagree on how \mathbf{x} should be labeled. Getting its label allows us to get rid of either f or f' , so $VS_{\mathcal{F}}(S)$ and $DIS_{\mathcal{F}}(S)$ both shrink.

■ Recall that f^* is always compatible with examples (\mathbf{x}, y) , so it is always in $\mathcal{V} \rightarrow$ algorithm *zooms* into it!

■ This algorithm makes **no useless queries!** ... can we do better if we can choose \mathbf{x} ?

Input: models \mathcal{F}

Output: selected model $f \in \mathcal{F}$

- 1: $L \leftarrow \emptyset$
- 2: $\mathcal{V} \leftarrow \mathcal{F}$ ▷ implements the version space $VS_{\mathcal{F}}(L)$
- 3: **for** $t = 1, 2, \dots, T$ **do**
- 4: $\mathbf{x} \leftarrow \operatorname{argmax}_{\mathbf{x} \in U} \operatorname{acq}_{VS}(\mathcal{V}, \mathcal{F}, \mathbf{x})$
- 5: obtain label y of \mathbf{x}
- 6: update $\mathcal{V} \leftarrow \{f \in \mathcal{V} : f(\mathbf{x}) = y\}$ ▷ update version space
- 7: **return** any $f \in \mathcal{V}$

■ We can always ensure that $\mathbf{x} \in DIS_{\mathcal{F}}(L)$ – as long as $VS_{\mathcal{F}}(S) \neq \emptyset$, in which case we can simply terminate.

Problem: how do we define the acquisition function?

Consider the *homogeneous linear classifiers*:

$$\mathcal{F} = \left\{ f_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbb{1}\{\boldsymbol{\theta}^\top \mathbf{x} > 0\} : \boldsymbol{\theta} \in \mathbb{R}^d, \|\boldsymbol{\theta}\|_2 = 1 \right\} \quad (28)$$

The version space of S is essentially the set of direction vectors $\boldsymbol{\theta}$ that classify all points correctly.

- Classifiers are hyperplanes in instance space and instances are hyperplanes in hypothesis space (**duality**)

■ Pick the point $x \in U$ that (greedily) restricts the version space as much as possible. In this special case, x it passes close to the **center** of $VS_{\mathcal{F}}(S)$.

Idea: pick $x \in U$ that reduces the *volume* of the version space $VS_{\mathcal{F}}(S)$ as much as possible.

Idea: pick $x \in U$ that reduces the *volume* of the version space $VS_{\mathcal{F}}(S)$ as much as possible.

■ The **volume** of a region $A \subseteq \mathcal{F}$ is:

$$\text{Vol}(A) = \int_A d\theta$$

Idea: pick $x \in U$ that reduces the *volume* of the version space $VS_{\mathcal{F}}(S)$ as much as possible.

■ The **volume** of a region $A \subseteq \mathcal{F}$ is:

$$\text{Vol}(A) = \int_A d\theta = \int_{\theta \in \mathbb{R}^{|\theta|}} \delta\{\theta \in A\} d\theta \quad (29)$$

So computing a volume in general requires **integration**.

Idea: pick $x \in U$ that reduces the *volume* of the version space $VS_{\mathcal{F}}(S)$ as much as possible.

■ The **volume** of a region $A \subseteq \mathcal{F}$ is:

$$\text{Vol}(A) = \int_A d\theta = \int_{\theta \in \mathbb{R}^{|\theta|}} \delta\{\theta \in A\} d\theta \quad (29)$$

So computing a volume in general requires **integration**.

■ Pick x that minimizes $\text{Vol}(VS_{\mathcal{F}}(S \cup \{(x, y)\}))$. However, we don't *know* the label y of x

Idea: pick $\mathbf{x} \in U$ that reduces the *volume* of the version space $VS_{\mathcal{F}}(S)$ as much as possible.

■ The **volume** of a region $A \subseteq \mathcal{F}$ is:

$$\text{Vol}(A) = \int_A d\boldsymbol{\theta} = \int_{\boldsymbol{\theta} \in \mathbb{R}^{|\boldsymbol{\theta}|}} \delta\{\boldsymbol{\theta} \in A\} d\boldsymbol{\theta} \quad (29)$$

So computing a volume in general requires **integration**.

■ Pick \mathbf{x} that minimizes $\text{Vol}(VS_{\mathcal{F}}(S \cup \{(\mathbf{x}, y)\}))$. However, we don't *know* the label y of $\mathbf{x} \rightarrow$ **marginalize it**:

$$\underset{\mathbf{x} \in U}{\text{argmin}} \mathbb{E}_{y \sim p_{\theta}(Y|\mathbf{x})} [\text{Vol}(VS_{\mathcal{F}}(S \cup \{(\mathbf{x}, y)\}))] \quad (30)$$

Question: how to **encode** the version space?

Question: how to **encode** the version space?

■ If \mathcal{F} is **finite**, can explicitly store $f \models S$. **Bonus:** computing expected volume is doable (integral becomes sum).

Question: how to **encode** the version space?

■ If \mathcal{F} is **finite**, can explicitly store $f \models S$. **Bonus:** computing expected volume is doable (integral becomes sum).

■ If \mathcal{F} is **infinite**, cannot store explicitly. However, we only need to compute its **volume**:

$$\mathbb{E}_{y \sim p_{\theta}(Y|x)} [\text{Vol}(\text{VS}_{\mathcal{F}}(S \cup \{(x, y)\}))] = \sum_{y \in [c]} p_{\theta}(y | x) \cdot \underbrace{\text{Vol}(\text{VS}_{\mathcal{F}}(S \cup \{(x, y)\}))}_{\text{this is the difficult bit}} \quad (31)$$

Question: how to **encode** the version space?

■ If \mathcal{F} is **finite**, can explicitly store $f \models S$. **Bonus:** computing expected volume is doable (integral becomes sum).

■ If \mathcal{F} is **infinite**, cannot store explicitly. However, we only need to compute its **volume**:

$$\mathbb{E}_{y \sim p_{\theta}(Y|x)} [\text{Vol}(VS_{\mathcal{F}}(S \cup \{(x, y)\}))] = \sum_{y \in [c]} p_{\theta}(y | x) \cdot \underbrace{\text{Vol}(VS_{\mathcal{F}}(S \cup \{(x, y)\}))}_{\text{this is the difficult bit}} \quad (31)$$

- If \mathcal{F} is **“simple”** and/or S is **small**, volume can be approximated cheaply using Monte Carlo techniques. For instance with *rejection sampling*, let $B \subseteq VS_{\mathcal{F}}(S)$ of known volume:

$$\{\tilde{\theta}_i \sim \text{Uniform}(B) : i \in [s]\}, \quad \text{Vol}(VS_{\mathcal{F}}(S')) \approx \frac{1}{\text{Vol}(B)} \cdot \frac{1}{s} \sum_{i \in [s]} \mathbb{1}\{\tilde{\theta}_i \in VS_{\mathcal{F}}(S')\} \quad (32)$$

To check, $\mathbb{1}\{\theta_i \in VS_{\mathcal{F}}(S')\}$, check that f_{θ} classifies all examples in S correctly.

Question: how to **encode** the version space?

■ If \mathcal{F} is **finite**, can explicitly store $f \models S$. **Bonus:** computing expected volume is doable (integral becomes sum).

■ If \mathcal{F} is **infinite**, cannot store explicitly. However, we only need to compute its **volume**:

$$\mathbb{E}_{y \sim p_{\theta}(Y|x)} [\text{Vol}(VS_{\mathcal{F}}(S \cup \{(x, y)\}))] = \sum_{y \in [c]} p_{\theta}(y | x) \cdot \underbrace{\text{Vol}(VS_{\mathcal{F}}(S \cup \{(x, y)\}))}_{\text{this is the difficult bit}} \quad (31)$$

- If \mathcal{F} is **“simple”** and/or S is **small**, volume can be approximated cheaply using Monte Carlo techniques. For instance with *rejection sampling*, let $B \subseteq VS_{\mathcal{F}}(S)$ of known volume:

$$\{\tilde{\theta}_i \sim \text{Uniform}(B) : i \in [s]\}, \quad \text{Vol}(VS_{\mathcal{F}}(S')) \approx \frac{1}{\text{Vol}(B)} \cdot \frac{1}{s} \sum_{i \in [s]} \mathbb{1}\{\tilde{\theta}_i \in VS_{\mathcal{F}}(S')\} \quad (32)$$

To check, $\mathbb{1}\{\theta_i \in VS_{\mathcal{F}}(S')\}$, check that f_{θ} classifies all examples in S correctly.

- Otherwise (think CNN on ImageNet), can be **extremely challenging**

■ These approaches make two **assumptions**:

■ These approaches make two **assumptions**:

- Disagreement is measured using *all* hypotheses in the version space $VS_{\mathcal{F}}(S)$.

■ These approaches make two **assumptions**:

- Disagreement is measured using *all* hypotheses in the version space $VS_{\mathcal{F}}(S)$.
- Disagreement is binary: it is only 0 if all hypotheses fully agree on $x \in U$.

■ Let's **relax** both of them \rightarrow **speed-up**!

■ These approaches make two **assumptions**:

- Disagreement is measured using *all* hypotheses in the version space $VS_{\mathcal{F}}(S)$.
- Disagreement is binary: it is only 0 if all hypotheses fully agree on $x \in U$.

■ Let's **relax** both of them \rightarrow **speed-up**!

■ Moreover, version space is only non-empty in the realizable case. How do we deal with this?

Idea:

- Subsample k representatives $C \subset VS_{\mathcal{F}}(S)$, $k \gg 1$
- Measure disagreement efficiently between representatives in C

Idea:

- Subsample k representatives $C \subset VS_{\mathcal{F}}(S)$, $k \gg 1$
- Measure disagreement efficiently between representatives in C

■ How to generate C ?

- Sample $C := \{c_j : j \in [k]\} \stackrel{k}{\sim} VS_{\mathcal{F}}(S)$, as done previously (i.e., we don't *gain* much)

Idea:

- Subsample k representatives $C \subset VS_{\mathcal{F}}(S)$, $k \gg 1$
- Measure disagreement efficiently between representatives in C

■ How to generate C ?

- Sample $\mathcal{C} := \{c_j : j \in [k]\} \stackrel{k}{\sim} VS_{\mathcal{F}}(S)$, as done previously (i.e., we don't *gain* much)
- Define a distribution $p(\theta | L)$, very Bayesian \rightarrow can be challenging to sample from the posterior

Idea:

- Subsample k representatives $C \subset VS_{\mathcal{F}}(S)$, $k \gg 1$
- Measure disagreement efficiently between representatives in C

■ How to generate C ?

- Sample $\mathcal{C} := \{c_j : j \in [k]\} \stackrel{k}{\sim} VS_{\mathcal{F}}(S)$, as done previously (i.e., we don't *gain* much)
- Define a distribution $p(\theta | L)$, very Bayesian \rightarrow can be challenging to sample from the posterior
- **Bagging**: *sample* k subsets of L , train one classifier c_j on each.

Idea:

- Subsample k representatives $C \subset VS_{\mathcal{F}}(S)$, $k \gg 1$
- Measure disagreement efficiently between representatives in C

■ How to generate C ?

- Sample $C := \{c_j : j \in [k]\} \stackrel{k}{\sim} VS_{\mathcal{F}}(S)$, as done previously (i.e., we don't *gain* much)
- Define a distribution $p(\theta | L)$, very Bayesian \rightarrow can be challenging to sample from the posterior
- **Bagging**: *sample* k subsets of L , train one classifier c_j on each.
- **Boosting**: *randomly reweight* L , sequentially train k classifiers by repeatedly reweighting examples by mistakes made by previous classifier.

■ How to measure disagreement of \mathcal{C} on $x \in U$?

■ How to measure disagreement of \mathcal{C} on $\mathbf{x} \in U$?

■ **“Hard” Voting + Entropy:**

$$\operatorname{argmax}_{\mathbf{x} \in U} - \sum_y \frac{n(y, \mathbf{x})}{k} \log \frac{n(y, \mathbf{x})}{k}, \quad n(y, \mathbf{x}) := \sum_{c \in \mathcal{C}} \mathbb{1}\{c(\mathbf{x}) = y\} \quad (33)$$

Each classifiers votes either 0 or 1.

■ How to measure disagreement of \mathcal{C} on $\mathbf{x} \in U$?

■ **“Hard” Voting + Entropy:**

$$\operatorname{argmax}_{\mathbf{x} \in U} - \sum_y \frac{n(y, \mathbf{x})}{k} \log \frac{n(y, \mathbf{x})}{k}, \quad n(y, \mathbf{x}) := \sum_{c \in \mathcal{C}} \mathbb{1}\{c(\mathbf{x}) = y\} \quad (33)$$

Each classifiers votes either 0 or 1.

■ **“Soft” Voting + Entropy:**

$$\operatorname{argmax}_{\mathbf{x} \in U} - \sum_y p_{\mathcal{C}}(y | \mathbf{x}) \log p_{\mathcal{C}}(y | \mathbf{x}), \quad p_{\mathcal{C}}(y | \mathbf{x}) := \frac{1}{k} \sum_{c \in \mathcal{C}} p_c(y | \mathbf{x}) \quad (34)$$

Output probabilities of each $c \in \mathcal{C}$ taken into account.

■ How to measure disagreement of \mathcal{C} on $\mathbf{x} \in U$?

■ “Hard” Voting + Entropy:

$$\operatorname{argmax}_{\mathbf{x} \in U} - \sum_y \frac{n(y, \mathbf{x})}{k} \log \frac{n(y, \mathbf{x})}{k}, \quad n(y, \mathbf{x}) := \sum_{c \in \mathcal{C}} \mathbb{1}\{c(\mathbf{x}) = y\} \quad (33)$$

Each classifiers votes either 0 or 1.

■ “Soft” Voting + Entropy:

$$\operatorname{argmax}_{\mathbf{x} \in U} - \sum_y p_{\mathcal{C}}(y | \mathbf{x}) \log p_{\mathcal{C}}(y | \mathbf{x}), \quad p_{\mathcal{C}}(y | \mathbf{x}) := \frac{1}{k} \sum_{c \in \mathcal{C}} p_c(y | \mathbf{x}) \quad (34)$$

Output probabilities of each $c \in \mathcal{C}$ taken into account.

■ Average Kullback-Liebler divergence:

$$\operatorname{argmax}_{\mathbf{x} \in U} \frac{1}{k} \sum_{c \in \mathcal{C}} \text{KL}(p_c(Y | \mathbf{x}) \| p_{\mathcal{C}}(Y | \mathbf{x})) \quad (35)$$

$$\text{KL}(p(Y | \mathbf{x}) \| q(Y | \mathbf{x})) := \sum_y p(y | \mathbf{x}) \log \frac{p(y | \mathbf{x})}{q(y | \mathbf{x})} \quad (36)$$

Very expressive, measures difference between **whole distributions**.

Model Improvement

Idea: pick the point that gives the maximal **improvement in model quality**

Model Improvement

Idea: pick the point that gives the maximal **improvement in model quality**

Useful Concepts

The **loss** of p_θ on example $z = (x, y)$ is denoted $\ell(\theta, z)$. For instance, cross-entropy loss:

$$\ell(\theta, z) := - \sum_j \mathbb{1}\{j = y\} \log p_\theta(j | x) = - \log p_\theta(y | x) \quad (37)$$

Idea: pick the point that gives the maximal **improvement in model quality**

Useful Concepts

The **loss** of p_θ on example $z = (x, y)$ is denoted $\ell(\theta, z)$. For instance, cross-entropy loss:

$$\ell(\theta, z) := - \sum_j \mathbb{1}\{j = y\} \log p_\theta(j | x) = - \log p_\theta(y | x) \quad (37)$$

The **true risk** \mathcal{L}^* of θ w.r.t. the *ground-truth distribution* $p^*(\mathbf{X}, Y)$ is:

$$\mathcal{L}^*(\theta) := \mathbb{E}_{z \sim p^*} [\ell(\theta, z)] \quad (38)$$

It measures the **true** quality of the model, unobserved.

Idea: pick the point that gives the maximal **improvement in model quality**

Useful Concepts

The **loss** of p_θ on example $z = (x, y)$ is denoted $\ell(\theta, z)$. For instance, cross-entropy loss:

$$\ell(\theta, z) := - \sum_j \mathbb{1}\{j = y\} \log p_\theta(j | x) = - \log p_\theta(y | x) \quad (37)$$

The **true risk** \mathcal{L}^* of θ w.r.t. the *ground-truth distribution* $p^*(\mathbf{X}, Y)$ is:

$$\mathcal{L}^*(\theta) := \mathbb{E}_{z \sim p^*} [\ell(\theta, z)] \quad (38)$$

It measures the **true** quality of the model, unobserved.

The **empirical risk** $\hat{\mathcal{L}}_S$ of θ w.r.t. data set $S = \{z_1, \dots, z_m\}$ sampled i.i.d. from p^* is:

$$\hat{\mathcal{L}}_S(\theta) := \frac{1}{|S|} \sum_{z \in S} \ell(\theta, z) \quad (39)$$

It **estimates** the quality of the model from a sample S , optimized during training.

Let $\hat{\theta}$ be the parameters obtained by training on S and $\hat{\theta}^{+z}$ those obtained by training on $S \cup \{z\}$, i.e.,

$$\hat{\theta} := \operatorname{argmin}_{\theta} \hat{\mathcal{L}}_S(\theta) \quad \hat{\theta}^{+z} := \operatorname{argmin}_{\theta} \hat{\mathcal{L}}_{S \cup \{z\}}(\theta) \quad (40)$$

where optimization is possibly *approximate*, e.g., based on SGD.

Let $\hat{\theta}$ be the parameters obtained by training on S and $\hat{\theta}^{+z}$ those obtained by training on $S \cup \{z\}$, i.e.,

$$\hat{\theta} := \operatorname{argmin}_{\theta} \hat{\mathcal{L}}_S(\theta) \quad \hat{\theta}^{+z} := \operatorname{argmin}_{\theta} \hat{\mathcal{L}}_{S \cup \{z\}}(\theta) \quad (40)$$

where optimization is possibly *approximate*, e.g., based on SGD.

Model Improvement

The **model improvement** (MI) given by a new example $z \notin S$ is the decrease in true risk:

$$\operatorname{acq}(\mathbf{x}) := \mathcal{L}^*(\hat{\theta}) - \mathcal{L}^*(\hat{\theta}^{+z}) \quad (41)$$

The higher, the better \rightarrow pick the $\mathbf{x} \in U$ that **maximizes** the improvement.

Model Improvement as Greedy Optimization

■ MI amounts to solving:

$$\operatorname{argmax}_{x \in U} \mathcal{L}^*(\hat{\theta}) - \mathcal{L}^*(\hat{\theta}^{+z}) = \operatorname{argmin}_{x \in U} \mathcal{L}^*(\hat{\theta}^{+z}) \quad (42)$$

It is *guaranteed* to find the best next candidate!

⁵**Note:** MI is *greedy, not optimal!* Non-greedy alternatives are conceptually better, but they also computationally infeasible and for this reason they are ignored in the AL literature.

Model Improvement as Greedy Optimization

- MI amounts to solving:

$$\operatorname{argmax}_{x \in U} \mathcal{L}^*(\hat{\theta}) - \mathcal{L}^*(\hat{\theta}^{+z}) = \operatorname{argmin}_{x \in U} \mathcal{L}^*(\hat{\theta}^{+z}) \quad (42)$$

It is *guaranteed* to find the best next candidate!

- MI is essentially a **greedy strategy** for solving:⁵

$$\operatorname{argmin}_{S \subseteq U} \mathcal{L}^*(\hat{\theta}) \quad (43)$$

$$\text{s.t. } |S| \leq \text{query budget} \quad (44)$$

In this view, **AL is a subset optimization problem**, and MI solves it directly.

⁵**Note:** MI is *greedy*, *not* optimal! Non-greedy alternatives are conceptually better, but they also computationally infeasible and for this reason they are ignored in the AL literature.

Model Improvement as Greedy Optimization

- MI amounts to solving:

$$\operatorname{argmax}_{x \in U} \mathcal{L}^*(\hat{\theta}) - \mathcal{L}^*(\hat{\theta}^{+z}) = \operatorname{argmin}_{x \in U} \mathcal{L}^*(\hat{\theta}^{+z}) \quad (42)$$

It is *guaranteed* to find the best next candidate!

- MI is essentially a **greedy strategy** for solving:⁵

$$\operatorname{argmin}_{S \subseteq U} \mathcal{L}^*(\hat{\theta}) \quad (43)$$

$$\text{s.t. } |S| \leq \text{query budget} \quad (44)$$

In this view, **AL is a subset optimization problem**, and MI solves it directly.

- Compare this to *uncertainty sampling*, which is not as sound

⁵**Note:** MI is *greedy*, *not* optimal! Non-greedy alternatives are conceptually better, but they also computationally infeasible and for this reason they are ignored in the AL literature.

■ We want to solve:

$$\operatorname{argmin}_{x \in U} \mathcal{L}^*(\hat{\theta}^+ z) \quad (45)$$

⁶The unlabeled set U is ideally pretty large, so the approximation is reasonable.

■ We want to solve:

$$\operatorname{argmin}_{\mathbf{x} \in U} \mathcal{L}^*(\hat{\theta}^{+z}) \quad (45)$$

Problem: $\mathcal{L}^*(\cdot)$ is an **integral over $\mathbf{x}' \in \mathbb{R}^d$** :

$$\mathcal{L}^*(\hat{\theta}^{+z}) = \mathbb{E}_{z' \sim p^*}[\ell(\hat{\theta}^{+z}, z')] = \int_{\mathbb{R}^d} \ell(\hat{\theta}^{+z}, (\mathbf{x}', y')) d\mathbf{x}' \quad (46)$$

which is intractable

⁶The unlabeled set U is ideally pretty large, so the approximation is reasonable.

■ We want to solve:

$$\operatorname{argmin}_{\mathbf{x} \in U} \mathcal{L}^*(\hat{\theta}^{+z}) \quad (45)$$

Problem: $\mathcal{L}^*(\cdot)$ is an **integral over $\mathbf{x}' \in \mathbb{R}^d$** :

$$\mathcal{L}^*(\hat{\theta}^{+z}) = \mathbb{E}_{z' \sim p^*} [\ell(\hat{\theta}^{+z}, z')] = \int_{\mathbb{R}^d} \ell(\hat{\theta}^{+z}, (\mathbf{x}', y')) d\mathbf{x}' \quad (46)$$

which is intractable \rightarrow **approximate using empirical average over U** :⁶

$$\mathcal{L}^*(\hat{\theta}^{+z}) \approx \hat{\mathcal{L}}_U(\hat{\theta}^{+z}) = \frac{1}{|U|} \sum_{\mathbf{x}' \in U} \ell(\hat{\theta}^{+z}, (\mathbf{x}', y')) \quad (47)$$

⁶The unlabeled set U is ideally pretty large, so the approximation is reasonable.

■ We want to solve:

$$\operatorname{argmin}_{\mathbf{x} \in U} \mathcal{L}^*(\hat{\theta}^{+z}) \quad (45)$$

Problem: $\mathcal{L}^*(\cdot)$ is an **integral over $\mathbf{x}' \in \mathbb{R}^d$** :

$$\mathcal{L}^*(\hat{\theta}^{+z}) = \mathbb{E}_{z' \sim p^*} [\ell(\hat{\theta}^{+z}, z')] = \int_{\mathbb{R}^d} \ell(\hat{\theta}^{+z}, (\mathbf{x}', y')) d\mathbf{x}' \quad (46)$$

which is intractable \rightarrow **approximate using empirical average over U** :⁶

$$\mathcal{L}^*(\hat{\theta}^{+z}) \approx \hat{\mathcal{L}}_U(\hat{\theta}^{+z}) = \frac{1}{|U|} \sum_{\mathbf{x}' \in U} \ell(\hat{\theta}^{+z}, (\mathbf{x}', y')) \quad (47)$$

Example: if ℓ is the 0–1 loss, then this amounts to:

$$\frac{1}{|U|} \sum_{\mathbf{x}' \in U} \mathbb{1}\{f_{\hat{\theta}^{+z}}(\mathbf{x}') \neq y'\} \quad (48)$$

⁶The unlabeled set U is ideally pretty large, so the approximation is reasonable.

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \hat{\mathcal{L}}_U(\hat{\theta}^{+z}) = \frac{1}{|U|} \sum_{x' \in U} \ell(\hat{\theta}^{+z}, (x', y')) \quad (49)$$

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \hat{\mathcal{L}}_U(\hat{\theta}^{+z}) = \frac{1}{|U|} \sum_{x' \in U} \ell(\hat{\theta}^{+z}, (x', y')) \quad (49)$$

Problem: we don't have access to the ground-truth label $z = (x, y)$

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \hat{\mathcal{L}}_U(\hat{\theta}^{+z}) = \frac{1}{|U|} \sum_{\mathbf{x}' \in U} \ell(\hat{\theta}^{+z}, (\mathbf{x}', \mathbf{y}')) \quad (49)$$

Problem: we don't have access to the ground-truth label $z = (\mathbf{x}, \mathbf{y}) \rightarrow$ **marginalize w.r.t. p^* :**

$$\mathbb{E}_{\mathbf{y} \sim p^*(Y|\mathbf{x})} \left[\frac{1}{|U|} \sum_{\mathbf{x}' \in U} \ell(\hat{\theta}^{+(\mathbf{x}, \mathbf{y})}, (\mathbf{x}', \mathbf{y}')) \right] \quad (50)$$

This averages over alternative **future models** $\hat{\theta}^{+(\mathbf{x}, \mathbf{y})}$ obtained after retraining on $L \cup (\mathbf{x}, \mathbf{y})$.

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \hat{\mathcal{L}}_U(\hat{\theta}^{+z}) = \frac{1}{|U|} \sum_{\mathbf{x}' \in U} \ell(\hat{\theta}^{+z}, (\mathbf{x}', \mathbf{y}')) \quad (49)$$

Problem: we don't have access to the ground-truth label $z = (\mathbf{x}, \mathbf{y}) \rightarrow$ **marginalize w.r.t. p^* :**

$$\mathbb{E}_{\mathbf{y} \sim p^*(Y|\mathbf{x})} \left[\frac{1}{|U|} \sum_{\mathbf{x}' \in U} \ell(\hat{\theta}^{+(\mathbf{x}, \mathbf{y})}, (\mathbf{x}', \mathbf{y}')) \right] \quad (50)$$

This averages over alternative **future models** $\hat{\theta}^{+(\mathbf{x}, \mathbf{y})}$ obtained after retraining on $L \cup (\mathbf{x}, \mathbf{y})$.

Problem: we don't have access to the ground-truth label \mathbf{y}' either

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \hat{\mathcal{L}}_U(\hat{\theta}^{+z}) = \frac{1}{|U|} \sum_{x' \in U} \ell(\hat{\theta}^{+z}, (x', y')) \quad (49)$$

Problem: we don't have access to the ground-truth label $z = (x, y) \rightarrow$ **marginalize w.r.t. p^* :**

$$\mathbb{E}_{y \sim p^*(Y|x)} \left[\frac{1}{|U|} \sum_{x' \in U} \ell(\hat{\theta}^{+(x,y)}, (x', y')) \right] \quad (50)$$

This averages over alternative **future models** $\hat{\theta}^{+(x,y)}$ obtained after retraining on $L \cup (x, y)$.

Problem: we don't have access to the ground-truth label y' either \rightarrow **marginalize w.r.t. p^* :**

$$\mathbb{E}_{y \sim p^*(Y|x)} \left[\frac{1}{|U|} \sum_{x' \in U} \mathbb{E}_{y' \sim p^*(Y|x')} \left[\ell(\hat{\theta}^{+(x,y)}, (x', y')) \right] \right] \quad (51)$$

This averages over the **unknown labels** y' of the instances in $x' \in U$.

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p^*(Y|\mathbf{x})} \left[\frac{1}{|U|} \sum_{\mathbf{x}' \in U} \mathbb{E}_{y' \sim p^*(Y|\mathbf{x}')} \left[\ell(\hat{\theta}^{+(x,y)}, (\mathbf{x}', y')) \right] \right] \quad (52)$$

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p^*(Y|\mathbf{x})} \left[\frac{1}{|U|} \sum_{\mathbf{x}' \in U} \mathbb{E}_{y' \sim p^*(Y|\mathbf{x}')} \left[\ell(\hat{\theta}^{+(x,y)}, (\mathbf{x}', y')) \right] \right] \quad (52)$$

Problem: we don't have access to p^* at all

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \rightarrow \quad \mathbb{E}_{y \sim p^*(Y|x)} \left[\frac{1}{|U|} \sum_{x' \in U} \mathbb{E}_{y' \sim p^*(Y|x')} \left[\ell(\hat{\theta}^{+(x,y)}, (x', y')) \right] \right] \quad (52)$$

Problem: we don't have access to p^* at all \rightarrow **estimate using model's distribution:**

$$\mathbb{E}_{y \sim p_{\hat{\theta}}(Y|x)} \left[\underbrace{\frac{1}{|U|} \sum_{x' \in U} \underbrace{\mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|x')} \left[\ell(\hat{\theta}^+, (x', y')) \right]}_{(a)}}_{(b)} \right] \quad (53)$$

(c)

where $\hat{\theta}^+ := \hat{\theta}^{+(x,y)}$. If $p_{\hat{\theta}}$ is "good enough", then the approximation is good.

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p^*(Y|x)} \left[\frac{1}{|U|} \sum_{x' \in U} \mathbb{E}_{y' \sim p^*(Y|x')} \left[\ell(\hat{\theta}^{+(x,y)}, (x', y')) \right] \right] \quad (52)$$

Problem: we don't have access to p^* at all \rightarrow **estimate using model's distribution:**

$$\mathbb{E}_{y \sim p_{\hat{\theta}}(Y|x)} \left[\underbrace{\frac{1}{|U|} \sum_{x' \in U} \underbrace{\mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|x')} \left[\ell(\hat{\theta}^+, (x', y')) \right]}_{(a)}}_{(b)} \right] \quad (53)$$

(c)

where $\hat{\theta}^+ := \hat{\theta}^{+(x,y)}$. If $p_{\hat{\theta}}$ is "good enough", then the approximation is good.

(a) Is the **expected** loss of the updated model on $x' \in U$,

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p^*(Y|x)} \left[\frac{1}{|U|} \sum_{x' \in U} \mathbb{E}_{y' \sim p^*(Y|x')} \left[\ell(\hat{\theta}^{+(x,y)}, (x', y')) \right] \right] \quad (52)$$

Problem: we don't have access to p^* at all \rightarrow **estimate using model's distribution:**

$$\mathbb{E}_{y \sim p_{\hat{\theta}}(Y|x)} \left[\underbrace{\frac{1}{|U|} \sum_{x' \in U} \underbrace{\mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|x')} \left[\ell(\hat{\theta}^+, (x', y')) \right]}_{(a)}}_{(b)} \right] \quad (53)$$

(c)

where $\hat{\theta}^+ := \hat{\theta}^{+(x,y)}$. If $p_{\hat{\theta}}$ is "good enough", then the approximation is good.

- (a) Is the **expected** loss of the updated model on $x' \in U$,
- (b) Is the **average** expected loss of the updated model on *all* of U ,

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p^*(Y|x)} \left[\frac{1}{|U|} \sum_{x' \in U} \mathbb{E}_{y' \sim p^*(Y|x')} \left[\ell(\hat{\theta}^{+(x,y)}, (x', y')) \right] \right] \quad (52)$$

Problem: we don't have access to p^* at all \rightarrow **estimate using model's distribution:**

$$\mathbb{E}_{y \sim p_{\hat{\theta}}(Y|x)} \left[\underbrace{\frac{1}{|U|} \sum_{x' \in U} \underbrace{\mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|x')} \left[\ell(\hat{\theta}^+, (x', y')) \right]}_{(a)}}_{(b)} \right] \quad (53)$$

(c)

where $\hat{\theta}^+ := \hat{\theta}^{+(x,y)}$. If $p_{\hat{\theta}}$ is “good enough”, then the approximation is good.

- (a) Is the **expected** loss of the updated model on $x' \in U$,
- (b) Is the **average** expected loss of the updated model on *all* of U ,
- (c) Is the above **averaged over the possible updated models** $\hat{\theta}^{+(x,y)}$.

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p_{\hat{\theta}}(Y|x)} \left[\frac{1}{|U|} \sum_{x' \in U} \mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|x')} \left[\ell(\hat{\theta}^+, (x', y')) \right] \right] \quad (54)$$

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p_{\hat{\theta}}(Y|\mathbf{x})} \left[\frac{1}{|U|} \sum_{\mathbf{x}' \in U} \mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|\mathbf{x}')} \left[\ell(\hat{\theta}^+, (\mathbf{x}', y')) \right] \right] \quad (54)$$

Example: consider the 0-1 loss $\ell(\theta, (\mathbf{x}, y)) = \mathbb{1}\{f_{\theta}(\mathbf{x}) \neq y\}$.

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p_{\hat{\theta}}(Y|\mathbf{x})} \left[\frac{1}{|U|} \sum_{\mathbf{x}' \in U} \mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|\mathbf{x}')} \left[\ell(\hat{\theta}^+, (\mathbf{x}', y')) \right] \right] \quad (54)$$

Example: consider the 0-1 loss $\ell(\theta, (\mathbf{x}, y)) = \mathbb{1}\{f_{\theta}(\mathbf{x}) \neq y\}$. Then:

$$\mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|\mathbf{x}')} \left[\mathbb{1}\{f_{\hat{\theta}^+}(\mathbf{x}') \neq y'\} \right] = p_{\hat{\theta}^+}(\hat{y}' \neq y' | \mathbf{x}'), \quad \hat{y}' := f_{\hat{\theta}^+}(\mathbf{x}') \quad (55)$$

$$= 1 - p_{\hat{\theta}^+}(\hat{y}' | \mathbf{x}') \quad (56)$$

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p_{\hat{\theta}}(Y|\mathbf{x})} \left[\frac{1}{|U|} \sum_{\mathbf{x}' \in U} \mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|\mathbf{x}')} \left[\ell(\hat{\theta}^+, (\mathbf{x}', y')) \right] \right] \quad (54)$$

Example: consider the 0-1 loss $\ell(\theta, (\mathbf{x}, y)) = \mathbb{1}\{f_{\theta}(\mathbf{x}) \neq y\}$. Then:

$$\mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|\mathbf{x}')} \left[\mathbb{1}\{f_{\hat{\theta}^+}(\mathbf{x}') \neq y'\} \right] = p_{\hat{\theta}^+}(\hat{y}' \neq y' | \mathbf{x}'), \quad \hat{y}' := f_{\hat{\theta}^+}(\mathbf{x}') \quad (55)$$

$$= 1 - p_{\hat{\theta}^+}(\hat{y}' | \mathbf{x}') \quad (56)$$

Hence, the above can be rewritten as ($\frac{1}{|U|}$ doesn't matter because it is independent of \mathbf{x}):

$$\mathbb{E}_{y \sim p_{\hat{\theta}}(Y|\mathbf{x})} \left[\frac{1}{|U|} \sum_{\mathbf{x}' \in U} (1 - p_{\hat{\theta}^+}(\hat{y}' | \mathbf{x}')) \right] \propto \sum_{y \in [c]} p_{\hat{\theta}}(y | \mathbf{x}) \sum_{\mathbf{x}' \in U} (1 - p_{\hat{\theta}^+}(\hat{y}' | \mathbf{x}')) \quad (57)$$

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p_{\hat{\theta}}(Y|\mathbf{x})} \left[\frac{1}{|U|} \sum_{\mathbf{x}' \in U} \mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|\mathbf{x}')} \left[\ell(\hat{\theta}^+, (\mathbf{x}', y')) \right] \right] \quad (54)$$

Example: consider the 0-1 loss $\ell(\theta, (\mathbf{x}, y)) = \mathbb{1}\{f_{\theta}(\mathbf{x}) \neq y\}$. Then:

$$\mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|\mathbf{x}')} \left[\mathbb{1}\{f_{\hat{\theta}^+}(\mathbf{x}') \neq y'\} \right] = p_{\hat{\theta}^+}(\hat{y}' \neq y' | \mathbf{x}'), \quad \hat{y}' := f_{\hat{\theta}^+}(\mathbf{x}') \quad (55)$$

$$= 1 - p_{\hat{\theta}^+}(\hat{y}' | \mathbf{x}') \quad (56)$$

Hence, the above can be rewritten as ($\frac{1}{|U|}$ doesn't matter because it is independent of \mathbf{x}):

$$\mathbb{E}_{y \sim p_{\hat{\theta}}(Y|\mathbf{x})} \left[\frac{1}{|U|} \sum_{\mathbf{x}' \in U} (1 - p_{\hat{\theta}^+}(\hat{y}' | \mathbf{x}')) \right] \propto \sum_{y \in [c]} p_{\hat{\theta}}(y | \mathbf{x}) \sum_{\mathbf{x}' \in U} (1 - p_{\hat{\theta}^+}(\hat{y}' | \mathbf{x}')) \quad (57)$$

■ We pick $\mathbf{x} \in U$ that minimizes the above \rightarrow **minimizes expected future confidence on U**

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p_{\hat{\theta}}(Y|x)} \left[\frac{1}{|U|} \sum_{x' \in U} \mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|x')} \left[\ell(\hat{\theta}^+, (x', y')) \right] \right] \quad (58)$$

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p_{\hat{\theta}}(Y|\mathbf{x})} \left[\frac{1}{|U|} \sum_{\mathbf{x}' \in U} \mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|\mathbf{x}')} \left[\ell(\hat{\theta}^+, (\mathbf{x}', y')) \right] \right] \quad (58)$$

Example: consider the negative log-likelihood $\ell(\theta, (\mathbf{x}, y)) = -\log p_{\theta}(y | \mathbf{x})$.

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p_{\hat{\theta}}(Y|x)} \left[\frac{1}{|U|} \sum_{x' \in U} \mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|x')} \left[\ell(\hat{\theta}^+, (x', y')) \right] \right] \quad (58)$$

Example: consider the negative log-likelihood $\ell(\theta, (x, y)) = -\log p_{\theta}(y | x)$. Then:

$$\mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|x')} \left[-\log p_{\hat{\theta}^+}(y' | x') \right] = -\sum_{y' \in [c]} p_{\hat{\theta}^+}(y' | x') \log p_{\hat{\theta}^+}(y' | x') \quad (59)$$

$$= H_{\hat{\theta}^+}(Y | x) \quad (60)$$

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p_{\hat{\theta}}(Y|x)} \left[\frac{1}{|U|} \sum_{x' \in U} \mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|x')} \left[\ell(\hat{\theta}^+, (x', y')) \right] \right] \quad (58)$$

Example: consider the negative log-likelihood $\ell(\theta, (x, y)) = -\log p_{\theta}(y | x)$. Then:

$$\mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|x')} \left[-\log p_{\hat{\theta}^+}(y' | x') \right] = -\sum_{y' \in [c]} p_{\hat{\theta}^+}(y' | x') \log p_{\hat{\theta}^+}(y' | x') \quad (59)$$

$$= H_{\hat{\theta}^+}(Y | x) \quad (60)$$

Hence, the above can be rewritten as:

$$\mathbb{E}_{y \sim p_{\hat{\theta}}(Y|x)} \left[\frac{1}{|U|} \sum_{x' \in U} (H_{\hat{\theta}^+}(Y | x)) \right] \propto \sum_{y \in [c]} p_{\hat{\theta}}(y | x) \sum_{x' \in U} (H_{\hat{\theta}^+}(Y | x)) \quad (61)$$

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p_{\hat{\theta}}(Y|x)} \left[\frac{1}{|U|} \sum_{x' \in U} \mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|x')} \left[\ell(\hat{\theta}^+, (x', y')) \right] \right] \quad (58)$$

Example: consider the negative log-likelihood $\ell(\theta, (x, y)) = -\log p_{\theta}(y | x)$. Then:

$$\mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|x')} \left[-\log p_{\hat{\theta}^+}(y' | x') \right] = -\sum_{y' \in [c]} p_{\hat{\theta}^+}(y' | x') \log p_{\hat{\theta}^+}(y' | x') \quad (59)$$

$$= H_{\hat{\theta}^+}(Y | x) \quad (60)$$

Hence, the above can be rewritten as:

$$\mathbb{E}_{y \sim p_{\hat{\theta}}(Y|x)} \left[\frac{1}{|U|} \sum_{x' \in U} (H_{\hat{\theta}^+}(Y | x)) \right] \propto \sum_{y \in [c]} p_{\hat{\theta}}(y | x) \sum_{x' \in U} (H_{\hat{\theta}^+}(Y | x)) \quad (61)$$

■ We pick $x \in U$ that minimizes the above \rightarrow **minimizes expected future entropy on U**

- In uncertainty sampling, we pick x that minimizes model's estimate of current uncertainty w.r.t. itself, this is **myopic**
- In expected model improvement, we pick x that minimizes model's estimate of expected future uncertainty w.r.t. unlabeled set, this is **less myopic**

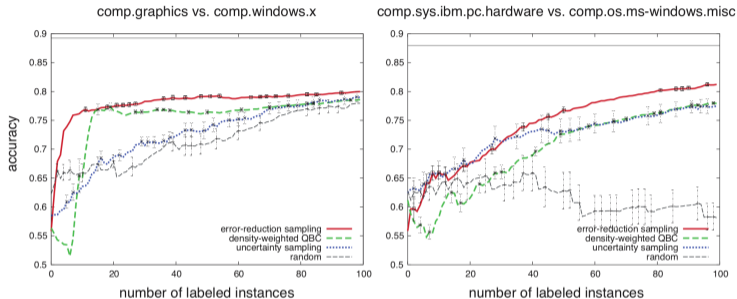


Figure 4.1: Learning curves showing that expected error reduction can outperform QBC and uncertainty sampling for two binary text classification tasks. *Source:* Adapted from Roy and McCallum (2001), with kind permission of the authors.

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p_{\hat{\theta}}(Y|x)} \left[\frac{1}{|U|} \sum_{x' \in U} \mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|x')} \left[\ell(\hat{\theta}^+, (x', y')) \right] \right] \quad (62)$$

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p_{\hat{\theta}}(Y|x)} \left[\frac{1}{|U|} \sum_{x' \in U} \mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|x')} \left[\ell(\hat{\theta}^+, (x', y')) \right] \right] \quad (62)$$

Problem: computing $\hat{\theta}^+$ requires to **fit model on** $L \cup \{(x, y)\}$ (**slow**)

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p_{\hat{\theta}}(Y|x)} \left[\frac{1}{|U|} \sum_{x' \in U} \mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|x')} \left[\ell(\hat{\theta}^+, (x', y')) \right] \right] \quad (62)$$

Problem: computing $\hat{\theta}^+$ requires to **fit model on** $L \cup \{(x, y)\}$ (**slow**)

Problem: this has to be done $|U| \times [c]$ times.

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p_{\hat{\theta}}(Y|x)} \left[\frac{1}{|U|} \sum_{x' \in U} \mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|x')} \left[\ell(\hat{\theta}^+, (x', y')) \right] \right] \quad (62)$$

Problem: computing $\hat{\theta}^+$ requires to **fit model on** $L \cup \{(x, y)\}$ (**slow**)

Problem: this has to be done $|U| \times [c]$ times.

Problem: this has to be done in each iteration of active learning.

■ We already decided on this approximation:

$$\mathcal{L}^*(\hat{\theta}^{+z}) \quad \longrightarrow \quad \mathbb{E}_{y \sim p_{\hat{\theta}}(Y|x)} \left[\frac{1}{|U|} \sum_{x' \in U} \mathbb{E}_{y' \sim p_{\hat{\theta}^+}(Y|x')} \left[\ell(\hat{\theta}^+, (x', y')) \right] \right] \quad (62)$$

Problem: computing $\hat{\theta}^+$ requires to **fit model on** $L \cup \{(x, y)\}$ (**slow**)

Problem: this has to be done $|U| \times [c]$ times.

Problem: this has to be done in each iteration of active learning.

■ Only **practical** for classes of models that support **closed-form updates** (e.g., Gaussian Processes) or stable **incremental learning** (e.g., perceptron-like learning algorithms).

Expected Model Change

- Unless a candidate (x, y) induces a large change in the model $\hat{\theta}$ upon retraining, then it cannot reduce the model's risk by much: **change is a prerequisite for improvement**.

Expected Model Change

■ Unless a candidate (x, y) induces a large change in the model $\hat{\theta}$ upon retraining, then it cannot reduce the model's risk by much: **change is a prerequisite for improvement**.

Intuition:

$$\ell(\hat{\theta}, z') - \ell(\hat{\theta}^{+z}, z') \leq |\ell(\hat{\theta}, z') - \ell(\hat{\theta}^{+z}, z')| \leq c \cdot \|\hat{\theta} - \hat{\theta}^{+z}\|, \quad c > 0 \quad (63)$$

where $\|\cdot\|$ is, e.g., the Euclidean norm. This formally holds for all c -Lipshitz loss functions ℓ .

Expected Model Change

■ Unless a candidate (\mathbf{x}, \mathbf{y}) induces a large change in the model $\hat{\theta}$ upon retraining, then it cannot reduce the model's risk by much: **change is a prerequisite for improvement**.

Intuition:

$$\ell(\hat{\theta}, z') - \ell(\hat{\theta}^{+z}, z') \leq |\ell(\hat{\theta}, z') - \ell(\hat{\theta}^{+z}, z')| \leq c \cdot \|\hat{\theta} - \hat{\theta}^{+z}\|, \quad c > 0 \quad (63)$$

where $\|\cdot\|$ is, e.g., the Euclidean norm. This formally holds for all c -Lipshitz loss functions ℓ .

■ Large change also occurs when the loss *increases* – hence the absolute value in the second step of Eq. 63.

All in all, EMC looks for examples $\mathbf{x} \in U$ that “make a difference” one way or the other.

But once (\mathbf{x}, \mathbf{y}) is acquired it is added to the training set L on which $\hat{\theta}$ is fit, so loss is likely to go *down* rather than *up*.

■ The trick is that if $\hat{\theta}$ is obtained via **gradient descent**, the difference $\hat{\theta} - \hat{\theta}^{+z}$ is easy to compute:

$$\hat{\theta} - \hat{\theta}^{+z} = \eta \cdot \nabla_{\theta} \ell(\theta, z) \quad (64)$$

where η is the learning rate. This gives **expected gradient length**:

$$\text{acc}_{\text{EGL}}(\mathbf{x}) := \mathbb{E}_{y \sim p_{\theta}(Y|\mathbf{x})} \left[\|\nabla_{\theta} \ell(\hat{\theta}, (\mathbf{x}, y))\|^2 \right] \quad (65)$$

The square does not change ranking of examples & avoids computing a square root.

- Quite cheap to compute using automatic differentiation packages (using Jacobian to parallelize over U)
- Assuming η is constant across examples and GD, the computation is exact. For other optimizers, it is an approximation

Are Uncertain Points Representative?

Idea: pick instances $\mathbf{x} \in U$ that are *both* **locally informative** and also **similar** to as many other unlabeled points as possible:

$$\operatorname{argmax}_{\mathbf{x} \in U} \operatorname{acq}(f, \mathbf{x}) \cdot \left(\frac{1}{|U|} \sum_{\mathbf{x}' \in U} \operatorname{sim}(\mathbf{x}, \mathbf{x}') \right)^{\beta} \quad (66)$$

where:

- $\operatorname{acq}(f, \mathbf{x})$ is a “standard” acquisition function based on, e.g., pointwise uncertainty.
- $\operatorname{sim}(\mathbf{x}, \mathbf{x}')$ measures the **similarity** between \mathbf{x} and \mathbf{x}' , e.g., a Gaussian kernel, Pearson’s correlation coefficient, Spearman’s rank correlation. **Application specific.**
- $\beta > 0$ is a hyper-parameter

Intuitively, \mathbf{x} ’s label conveys information about the label on the other points in U

Example

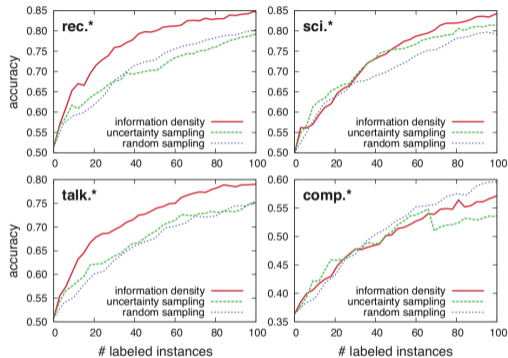


Figure 5.2: Learning curves showing that, by explicitly weighting queries by their representativeness among the input instances, information density can yield better results than the base uncertainty sampling heuristic by itself.

■ We optimize:

$$\operatorname{argmax}_{\mathbf{x} \in U} \operatorname{acq}(f, \mathbf{x}) \cdot \left(\frac{1}{|U|} \sum_{\mathbf{x}' \in U} \operatorname{sim}(\mathbf{x}, \mathbf{x}') \right)^\beta \quad (67)$$

Properties:

- Tends to work better than pure more “local” acquisition functions [Settles, 2012]
- Even when uncertainty sampling is worse than random, information density performs well

■ We optimize:

$$\operatorname{argmax}_{\mathbf{x} \in U} \operatorname{acq}(f, \mathbf{x}) \cdot \left(\frac{1}{|U|} \sum_{\mathbf{x}' \in U} \operatorname{sim}(\mathbf{x}, \mathbf{x}') \right)^\beta \quad (67)$$

Properties:

- Tends to work better than pure more “local” acquisition functions [Settles, 2012]
- Even when uncertainty sampling is worse than random, information density performs well
- Similarity computation can be sped-up using **caching**: “simply” store similarity matrix $S_{ij} = [\operatorname{sim}(\mathbf{x}_i, \mathbf{x}_j)]$ for all $\mathbf{x}_i, \mathbf{x}_j \in U$ (only needs to be done **once**)

■ We optimize:

$$\operatorname{argmax}_{\mathbf{x} \in U} \operatorname{acq}(f, \mathbf{x}) \cdot \left(\frac{1}{|U|} \sum_{\mathbf{x}' \in U} \operatorname{sim}(\mathbf{x}, \mathbf{x}') \right)^\beta \quad (67)$$

Properties:

- Tends to work better than pure more “local” acquisition functions [Settles, 2012]
- Even when uncertainty sampling is worse than random, information density performs well
- Similarity computation can be sped-up using **caching**: “simply” store similarity matrix $S_{ij} = [\operatorname{sim}(\mathbf{x}_i, \mathbf{x}_j)]$ for all $\mathbf{x}_i, \mathbf{x}_j \in U$ (only needs to be done **once**)
- Approximate using **clustering**: cluster U so that points within cluster are similar and points across clusters are not \rightarrow block-diagonal similarity matrix, lowers storage requirement from $O(|U|^2)$ to $O(\sum_i |\text{cluster}_i|^2)$

■ Do we gain anything by “summarizing” the data using **clustering**?

■ Do we gain anything by “summarizing” the data using **clustering**?

Idea:

- Cluster unlabeled data set $U \rightarrow \{C_i \subset U : i \in [k]\}$
- Treat each C_i as a separate problem, e.g., **query cluster centroids**

- Do we gain anything by “summarizing” the data using **clustering**?

Idea:

- Cluster unlabeled data set $U \rightarrow \{C_i \subset U : i \in [k]\}$
- Treat each C_i as a separate problem, e.g., **query cluster centroids**

Problems:

- U may not have a good clustering structure or $\text{sim}(\cdot, \cdot)$ may not be able to capture it
- How many clusters and at what granularity?
- Clusters of x 's may not correlate well with label y .

Figure: the swiss roll dataset has no obvious clustering structure.

WRITE ME

Extensions

■ Consider a neural network $f_\theta : \mathbb{R}^d \rightarrow [c]$:

$$f_\theta(\mathbf{x}) = \operatorname{argmax}_{y \in [c]} p_\theta(y | \mathbf{x})$$

$$p_\theta(y | \mathbf{x}) = \operatorname{softmax}(W\phi_\omega(\mathbf{x}))_y$$

where:

- $\theta = \{W, \omega\}$ are parameters
- $\phi_\omega : \mathbb{R}^d \rightarrow \mathbb{R}^k$ is an embedding function (e.g., convolutions + pooling layers)
- $W \in \mathbb{R}^{c \times k}$ are the parameters of the top dense layer

- Deep NNs have a number of **quirks**:

- Deep NNs have a number of **quirks**:
 - Very **overconfident** even away from the training set: their *uncertainty* cannot be trusted → strategies based on confidence, margin, entropy will underperform (including uncertainty sampling, model improvement, density-aware sampling, etc.)

- Deep NNs have a number of **quirks**:
 - Very **overconfident** even away from the training set: their *uncertainty* cannot be trusted → strategies based on confidence, margin, entropy will underperform (including uncertainty sampling, model improvement, density-aware sampling, etc.)
 - **Expensive** to fit on data: training a ResNet on a realistic data set can take minutes to days → hard to ensure responsiveness

- Deep NNs have a number of **quirks**:
 - Very **overconfident** even away from the training set: their *uncertainty* cannot be trusted → strategies based on confidence, margin, entropy will underperform (including uncertainty sampling, model improvement, density-aware sampling, etc.)
 - **Expensive** to fit on data: training a ResNet on a realistic data set can take minutes to days → hard to ensure responsiveness
 - Quite **insensitive** to the addition of a single example → what's the point of querying individual instances?

■ Deep NNs have a number of **quirks**:

- Very **overconfident** even away from the training set: their *uncertainty* cannot be trusted → strategies based on confidence, margin, entropy will underperform (including uncertainty sampling, model improvement, density-aware sampling, etc.)
- **Expensive** to fit on data: training a ResNet on a realistic data set can take minutes to days → hard to ensure responsiveness
- Quite **insensitive** to the addition of a single example → what's the point of querying individual instances?
- Training is **stochastic** (i.e., not 100% stable) → changes in performance can depend on factors other than new labeled examples, high variance

Overconfidence

Problem: Deep NNs tend to be very **overconfident** even away from the training set → their *uncertainty* cannot be trusted

Overconfidence

Problem: Deep NNs tend to be very **overconfident** even away from the training set → their *uncertainty* cannot be trusted

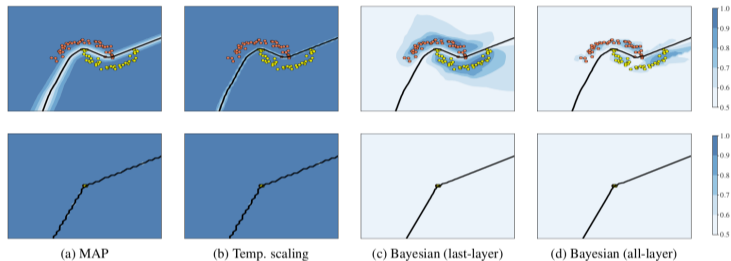


Figure 1. Binary classification on a toy dataset using a MAP estimate, temperature scaling, and both last-layer and all-layer Gaussian approximations over the weights which are obtained via Laplace approximations. Background color and black line represent confidence and decision boundary, respectively. Bottom row shows a zoomed-out view of the top row. The Bayesian approximations—even in the last-layer case—give desirable uncertainty estimates: confident close to the training data and uncertain otherwise. MAP and temperature scaling yield overconfident predictions. The optimal temperature is picked as in Guo et al. (2017).

Credit: [Kristiadi et al., 2020].

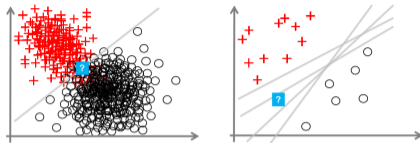


Figure 5: Left: Even with precise knowledge about the optimal hypothesis, the prediction at the query point (indicated by a question mark) is aleatorically uncertain, because the two classes are overlapping in that region. Right: A case of epistemic uncertainty due to a lack of knowledge about the right hypothesis, which is in turn caused by a lack of data.

Aleatoric vs Epistemic [Hüllermeier and Waegeman, 2021]

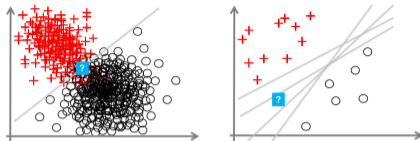


Figure 5: Left: Even with precise knowledge about the optimal hypothesis, the prediction at the query point (indicated by a question mark) is aleatorically uncertain, because the two classes are overlapping in that region. Right: A case of epistemic uncertainty due to a lack of knowledge about the right hypothesis, which is in turn caused by a lack of data.

■ **Aleatoric** uncertainty (“random”) captures how much we can trust the supervision itself. It cannot be decreased. (**left**)

Aleatoric vs Epistemic [Hüllermeier and Waegeman, 2021]

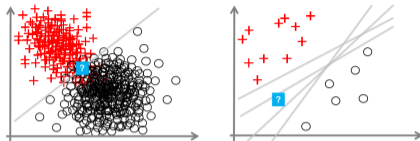


Figure 5: Left: Even with precise knowledge about the optimal hypothesis, the prediction at the query point (indicated by a question mark) is aleatorically uncertain, because the two classes are overlapping in that region. Right: A case of epistemic uncertainty due to a lack of knowledge about the right hypothesis, which is in turn caused by a lack of data.

■ **Aleatoric** uncertainty (“random”) captures how much we can trust the supervision itself. It cannot be decreased. (**left**)

■ **Epistemic** uncertainty (“relating to knowledge”) captures how little we know about the world. This reflects on uncertainty on the choice of θ . It decreases by acquiring **more data**. (**right**)

Aleatoric vs Epistemic [Hüllermeier and Waegeman, 2021]

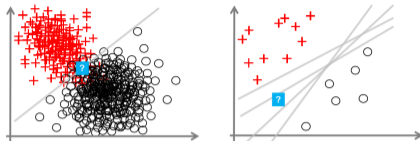


Figure 5: Left: Even with precise knowledge about the optimal hypothesis, the prediction at the query point (indicated by a question mark) is aleatorically uncertain, because the two classes are overlapping in that region. Right: A case of epistemic uncertainty due to a lack of knowledge about the right hypothesis, which is in turn caused by a lack of data.

■ **Aleatoric** uncertainty (“random”) captures how much we can trust the supervision itself. It cannot be decreased. (**left**)

■ **Epistemic** uncertainty (“relating to knowledge”) captures how little we know about the world. This reflects on uncertainty on the choice of θ . It decreases by acquiring **more data**. (**right**)

■ There isn't much point in trying to reduce **aleatoric** uncertainty in AL [Sharma and Bilgic, 2017]

Bayesian NNs

- The problem with NNs is that uncertainty depends on a *single model*:
 - This gives poor *epistemic* uncertainty

Bayesian NNs

- The problem with NNs is that uncertainty depends on a *single model*:
 - This gives poor *epistemic* uncertainty
 - Using **ensembles of NNs** is computationally challenging: training one NN is expensive, training k even more so

Bayesian NNs

- The problem with NNs is that uncertainty depends on a *single model*:
 - This gives poor *epistemic* uncertainty
 - Using **ensembles of NNs** is computationally challenging: training one NN is expensive, training k even more so
 - Using **Bayesian** techniques – i.e., maintaining a distribution over alternative NNs – is also challenging.

Bayesian NNs

- The problem with NNs is that uncertainty depends on a *single model*:
 - This gives poor *epistemic* uncertainty
 - Using **ensembles of NNs** is computationally challenging: training one NN is expensive, training k even more so
 - Using **Bayesian** techniques – i.e., maintaining a distribution over alternative NNs – is also challenging.

Idea of **Bayesian NNs**:

- Replace parameters θ with distribution over alternative parameters $p(\theta | L)$

Bayesian NNs

- The problem with NNs is that uncertainty depends on a *single model*:
 - This gives poor *epistemic* uncertainty
 - Using **ensembles of NNs** is computationally challenging: training one NN is expensive, training k even more so
 - Using **Bayesian** techniques – i.e., maintaining a distribution over alternative NNs – is also challenging.

Idea of **Bayesian NNs**:

- Replace parameters θ with distribution over alternative parameters $p(\theta | L)$
- Compute predictions by marginalizing over θ :

$$p(y | \mathbf{x}) = \int \underbrace{p(y | \mathbf{x}, \theta)}_{\text{NN with params } \theta} \cdot \underbrace{p(\theta | L)}_{\text{posterior over params}} d\theta \quad (68)$$

Bayesian NNs

- The problem with NNs is that uncertainty depends on a *single model*:
 - This gives poor *epistemic* uncertainty
 - Using **ensembles of NNs** is computationally challenging: training one NN is expensive, training k even more so
 - Using **Bayesian** techniques – i.e., maintaining a distribution over alternative NNs – is also challenging.

Idea of **Bayesian NNs**:

- Replace parameters θ with distribution over alternative parameters $p(\theta | L)$
- Compute predictions by marginalizing over θ :

$$p(y | \mathbf{x}) = \int \underbrace{p(y | \mathbf{x}, \theta)}_{\text{NN with params } \theta} \cdot \underbrace{p(\theta | L)}_{\text{posterior over params}} d\theta \quad (68)$$

- Learn by updating distribution:

$$p(\theta | L) \rightarrow p(\theta | L \cup \{(\mathbf{x}, y)\}) \quad (69)$$

Not trivial! Is there an **efficient approximation**?

Dropout

- Randomly set nodes to 0 with a fixed probability.

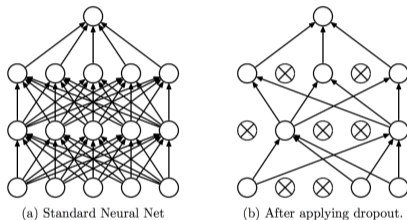


Figure 1: Dropout Neural Net Model. **Left:** A standard neural net with 2 hidden layers. **Right:** An example of a thinned net produced by applying dropout to the network on the left. Crossed units have been dropped.

- Used as a **regularization technique**: by randomly removing neurons, prevents them from relying on each other “too much”

Dropout as Bayesian Approximation

- Computing class probabilities:

$$p(y | \mathbf{x}, L) = \int p(y | \mathbf{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | L) d\boldsymbol{\theta} \quad (70)$$

$$\approx \int p(y | \mathbf{x}, \boldsymbol{\theta}) p_{\text{dropout}}(\boldsymbol{\theta}) d\boldsymbol{\theta} \quad (71)$$

$$\approx \frac{1}{R} \sum_{r=1}^R p(y | \mathbf{x}, \hat{\boldsymbol{\theta}}_r), \quad \hat{\boldsymbol{\theta}}_r \sim p_{\text{dropout}}(\boldsymbol{\theta}) \quad (72)$$

In other words, run NN R times with **dropout enabled** (during inference!) then average the R vectors of class probabilities.

Dropout as Bayesian Approximation

- Computing class probabilities:

$$p(y | \mathbf{x}, L) = \int p(y | \mathbf{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | L) d\boldsymbol{\theta} \quad (70)$$

$$\approx \int p(y | \mathbf{x}, \boldsymbol{\theta}) p_{\text{dropout}}(\boldsymbol{\theta}) d\boldsymbol{\theta} \quad (71)$$

$$\approx \frac{1}{R} \sum_{r=1}^R p(y | \mathbf{x}, \hat{\boldsymbol{\theta}}_r), \quad \hat{\boldsymbol{\theta}}_r \sim p_{\text{dropout}}(\boldsymbol{\theta}) \quad (72)$$

In other words, run NN R times with **dropout enabled** (during inference!) then average the R vectors of class probabilities.

- Dropout can be viewed as variational Bayesian approximation where the approximating distribution is a mixture of two Gaussians [Gal and Ghahramani, 2016].

Dropout as Bayesian Approximation

- Computing class probabilities:

$$p(y | \mathbf{x}, L) = \int p(y | \mathbf{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | L) d\boldsymbol{\theta} \quad (70)$$

$$\approx \int p(y | \mathbf{x}, \boldsymbol{\theta}) p_{\text{dropout}}(\boldsymbol{\theta}) d\boldsymbol{\theta} \quad (71)$$

$$\approx \frac{1}{R} \sum_{r=1}^R p(y | \mathbf{x}, \hat{\boldsymbol{\theta}}_r), \quad \hat{\boldsymbol{\theta}}_r \sim p_{\text{dropout}}(\boldsymbol{\theta}) \quad (72)$$

In other words, run NN R times with **dropout enabled** (during inference!) then average the R vectors of class probabilities.

- Dropout can be viewed as variational Bayesian approximation where the approximating distribution is a mixture of two Gaussians [Gal and Ghahramani, 2016]. The approximation is **independent** of $L \rightarrow$ no training required.

Dropout as Bayesian Approximation

- Computing class probabilities:

$$p(y | \mathbf{x}, L) = \int p(y | \mathbf{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | L) d\boldsymbol{\theta} \quad (70)$$

$$\approx \int p(y | \mathbf{x}, \boldsymbol{\theta}) p_{\text{dropout}}(\boldsymbol{\theta}) d\boldsymbol{\theta} \quad (71)$$

$$\approx \frac{1}{R} \sum_{r=1}^R p(y | \mathbf{x}, \hat{\boldsymbol{\theta}}_r), \quad \hat{\boldsymbol{\theta}}_r \sim p_{\text{dropout}}(\boldsymbol{\theta}) \quad (72)$$

In other words, run NN R times with **dropout enabled** (during inference!) then average the R vectors of class probabilities.

- Dropout can be viewed as variational Bayesian approximation where the approximating distribution is a mixture of two Gaussians [Gal and Ghahramani, 2016]. The approximation is **independent** of $L \rightarrow$ no training required.

- Immediately leads to **more calibrated output probabilities!**

Question: does dropout help with query selection too?

Question: does dropout help with query selection too? **Yes.**

■ Uncertainty sampling:

$$\text{acq}_{UNC}(\mathbf{x}) = - \sum_{y \in [c]} p(Y = y | \mathbf{x}, L) \log p(Y = y | \mathbf{x}, L) \quad (73)$$

Application is immediate.

Question: does dropout help with query selection too? **Yes.**

■ Uncertainty sampling:

$$\text{acq}_{UNC}(\mathbf{x}) = - \sum_{y \in [c]} p(Y = y | \mathbf{x}, L) \log p(Y = y | \mathbf{x}, L) \quad (73)$$

Application is immediate.

■ Mutual information between predictions and model posterior (**BALD**):

$$\text{acq}_{BALD}(\mathbf{x}) = H(Y | \mathbf{x}, L) - \mathbb{E}_{\theta \sim p(\theta | L)} [H(Y | \mathbf{x}, \theta)] \quad (74)$$

Question: does dropout help with query selection too? **Yes.**

■ Uncertainty sampling:

$$\text{acq}_{UNC}(\mathbf{x}) = - \sum_{y \in [c]} p(Y = y | \mathbf{x}, L) \log p(Y = y | \mathbf{x}, L) \quad (73)$$

Application is immediate.

■ Mutual information between predictions and model posterior (**BALD**):

$$\text{acq}_{BALD}(\mathbf{x}) = H(Y | \mathbf{x}, L) - \mathbb{E}_{\theta \sim p(\theta | L)} [H(Y | \mathbf{x}, \theta)] \quad (74)$$

- Left: entropy of the prediction \rightarrow high when the model's prediction is uncertain

Question: does dropout help with query selection too? **Yes.**

■ Uncertainty sampling:

$$\text{acq}_{UNC}(\mathbf{x}) = - \sum_{y \in [c]} p(Y = y | \mathbf{x}, L) \log p(Y = y | \mathbf{x}, L) \quad (73)$$

Application is immediate.

■ Mutual information between predictions and model posterior (**BALD**):

$$\text{acq}_{BALD}(\mathbf{x}) = H(Y | \mathbf{x}, L) - \mathbb{E}_{\theta \sim p(\theta | L)} [H(Y | \mathbf{x}, \theta)] \quad (74)$$

- Left: entropy of the prediction \rightarrow high when the model's prediction is uncertain
- Right: expected entropy of the prediction over the posterior of the model parameters \rightarrow low when the model is overall certain for each draw of model parameters from the posterior.

Question: does dropout help with query selection too? **Yes.**

■ Uncertainty sampling:

$$\text{acq}_{UNC}(\mathbf{x}) = - \sum_{y \in [c]} p(Y = y | \mathbf{x}, L) \log p(Y = y | \mathbf{x}, L) \quad (73)$$

Application is immediate.

■ Mutual information between predictions and model posterior (**BALD**):

$$\text{acq}_{BALD}(\mathbf{x}) = H(Y | \mathbf{x}, L) - \mathbb{E}_{\theta \sim p(\theta | L)} [H(Y | \mathbf{x}, \theta)] \quad (74)$$

- Left: entropy of the prediction \rightarrow high when the model's prediction is uncertain
- Right: expected entropy of the prediction over the posterior of the model parameters \rightarrow low when the model is overall certain for each draw of model parameters from the posterior.

High when model has many possible ways to explain the data, i.e., **the posterior draws are disagreeing among themselves.**

■ **BALD** can be computed as:

$$\begin{aligned}
 H(Y | \mathbf{x}, L) - \mathbb{E}_{\theta \sim p(\theta|L)}[H(Y|\mathbf{x}, \theta)] \\
 &= - \sum_{y \in [c]} p(y | \mathbf{x}, L) \log p(y | \mathbf{x}, L) + \mathbb{E}_{\theta \sim p(\theta|L)} \left[\sum_{y \in [c]} p(y | \mathbf{x}, \theta) \log p(y | \mathbf{x}, \theta) \right] \\
 &= - \sum_{y \in [c]} \int p(y | \mathbf{x}, \theta) p(\theta | L) d\theta \cdot \log \int p(y | \mathbf{x}, \theta) p(\theta | L) d\theta + \mathbb{E}_{\theta \sim p(\theta|L)} \left[\sum_{y \in [c]} p(y | \mathbf{x}, \theta) \log p(y | \mathbf{x}, \theta) \right] \\
 &\approx - \sum_{y \in [c]} \int p(y | \mathbf{x}, \theta) p_{drop}(\theta) d\theta \cdot \log \int p(y | \mathbf{x}, \theta) p_{drop}(\theta) d\theta + \mathbb{E}_{\theta \sim p_{drop}(\theta)} \left[\sum_{y \in [c]} p(y | \mathbf{x}, \theta) \log p(y | \mathbf{x}, \theta) \right] \\
 &\approx - \sum_{y \in [c]} \left(\frac{1}{R} \sum_r \hat{p}_y^r \right) \log \left(\frac{1}{R} \sum_r \hat{p}_y^r \right) + \frac{1}{R} \sum_{y,r} \hat{p}_y^r \log \hat{p}_y^r,
 \end{aligned}$$

where $\hat{p}_y^r = p(y | \mathbf{x}, \theta_r)$ is the output of the NN's softmax for class y and randomized parameters $\theta_r \sim p_{drop}(\theta)$.

■ **BALD** can be computed as:

$$\begin{aligned}
 & H(Y | \mathbf{x}, L) - \mathbb{E}_{\theta \sim p(\theta|L)}[H(Y|\mathbf{x}, \theta)] \\
 &= - \sum_{y \in [c]} p(y | \mathbf{x}, L) \log p(y | \mathbf{x}, L) + \mathbb{E}_{\theta \sim p(\theta|L)} \left[\sum_{y \in [c]} p(y | \mathbf{x}, \theta) \log p(y | \mathbf{x}, \theta) \right] \\
 &= - \sum_{y \in [c]} \int p(y | \mathbf{x}, \theta) p(\theta | L) d\theta \cdot \log \int p(y | \mathbf{x}, \theta) p(\theta | L) d\theta + \mathbb{E}_{\theta \sim p(\theta|L)} \left[\sum_{y \in [c]} p(y | \mathbf{x}, \theta) \log p(y | \mathbf{x}, \theta) \right] \\
 &\approx - \sum_{y \in [c]} \int p(y | \mathbf{x}, \theta) p_{drop}(\theta) d\theta \cdot \log \int p(y | \mathbf{x}, \theta) p_{drop}(\theta) d\theta + \mathbb{E}_{\theta \sim p_{drop}(\theta)} \left[\sum_{y \in [c]} p(y | \mathbf{x}, \theta) \log p(y | \mathbf{x}, \theta) \right] \\
 &\approx - \sum_{y \in [c]} \left(\frac{1}{R} \sum_r \hat{p}_y^r \right) \log \left(\frac{1}{R} \sum_r \hat{p}_y^r \right) + \frac{1}{R} \sum_{y,r} \hat{p}_y^r \log \hat{p}_y^r,
 \end{aligned}$$

where $\hat{p}_y^r = p(y | \mathbf{x}, \theta_r)$ is the output of the NN's softmax for class y and randomized parameters $\theta_r \sim p_{drop}(\theta)$.

■ Only need to compute \hat{p}^r once per $r \in [R]$, for a modest R times increase in runtime; could be parallelized.

■ **BALD** can be computed as:

$$\begin{aligned}
 & H(Y | \mathbf{x}, L) - \mathbb{E}_{\theta \sim p(\theta|L)}[H(Y|\mathbf{x}, \theta)] \\
 &= - \sum_{y \in [c]} p(y | \mathbf{x}, L) \log p(y | \mathbf{x}, L) + \mathbb{E}_{\theta \sim p(\theta|L)} \left[\sum_{y \in [c]} p(y | \mathbf{x}, \theta) \log p(y | \mathbf{x}, \theta) \right] \\
 &= - \sum_{y \in [c]} \int p(y | \mathbf{x}, \theta) p(\theta | L) d\theta \cdot \log \int p(y | \mathbf{x}, \theta) p(\theta | L) d\theta + \mathbb{E}_{\theta \sim p(\theta|L)} \left[\sum_{y \in [c]} p(y | \mathbf{x}, \theta) \log p(y | \mathbf{x}, \theta) \right] \\
 &\approx - \sum_{y \in [c]} \int p(y | \mathbf{x}, \theta) p_{drop}(\theta) d\theta \cdot \log \int p(y | \mathbf{x}, \theta) p_{drop}(\theta) d\theta + \mathbb{E}_{\theta \sim p_{drop}(\theta)} \left[\sum_{y \in [c]} p(y | \mathbf{x}, \theta) \log p(y | \mathbf{x}, \theta) \right] \\
 &\approx - \sum_{y \in [c]} \left(\frac{1}{R} \sum_r \hat{p}_y^r \right) \log \left(\frac{1}{R} \sum_r \hat{p}_y^r \right) + \frac{1}{R} \sum_{y,r} \hat{p}_y^r \log \hat{p}_y^r,
 \end{aligned}$$

where $\hat{p}_y^r = p(y | \mathbf{x}, \theta_r)$ is the output of the NN's softmax for class y and randomized parameters $\theta_r \sim p_{drop}(\theta)$.

■ Only need to compute \hat{p}^r once per $r \in [R]$, for a modest R times increase in runtime; could be parallelized.

■ As $R \rightarrow \infty$, this approximates the original mutual information.

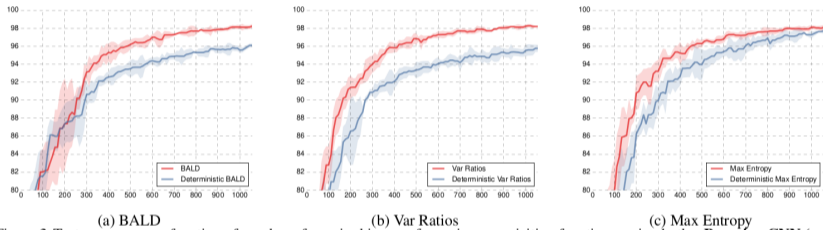


Figure 2. Test accuracy as a function of number of acquired images for various acquisition functions, using both a **Bayesian CNN (red)** and a **deterministic CNN (blue)**.

■ For all choices of acquisition function, the dropout-based uncertainty helps!

■ Let us look at **batch-based** active learning.

Batch Selection

Given L , U and a classifier $f \in \mathcal{F}$ trained on L , find a batch $B \subseteq U$ of $b \gg 1$ unlabeled instances that brings maximal information to the model:

$$\operatorname{argmax}_{B \subseteq U} \operatorname{acq}_{\text{BALD}}(f, B) \quad (75)$$

$$\text{s.t. } |B| = b \quad (76)$$

■ Let us look at **batch-based** active learning.

Batch Selection

Given L , U and a classifier $f \in \mathcal{F}$ trained on L , find a batch $B \subseteq U$ of $b \gg 1$ unlabeled instances that brings maximal information to the model:

$$\operatorname{argmax}_{B \subseteq U} \operatorname{acq}_{\text{BALD}}(f, B) \quad (75)$$

$$\text{s.t. } |B| = b \quad (76)$$

Advantages:

- Only retrain the model after ever b examples, meaning that supervision has an effect.
- Retraining is less frequent, leading to faster overall execution (at the expense of possibly instance selection, because b examples depend on a fixed f).
- Supports parallel annotation for, e.g., crowd-sourcing scenarios.

Question: can regular acquisition function (like BALD) be extended to this setting?

- Natural generalization of instance-level strategies:

$$\text{acq}(f, B) = \sum_{x \in B} \text{acq}(f, x) \quad (77)$$

How well does this work?

- Natural generalization of instance-level strategies:

$$\text{acq}(f, B) = \sum_{x \in B} \text{acq}(f, x) \quad (77)$$

How well does this work?

- This ignores *correlation* between instances in x :
 - Even if all of them are informative, they may **carry the same information**
 - We want B to be informative **as a whole!**

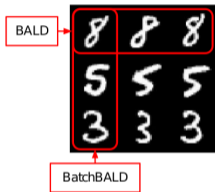


Figure 1: Idealised acquisitions of BALD and BatchBALD. If a dataset were to contain many (near) replicas for each data point, then BALD would select all replicas of a single informative data point at the expense of other informative data points, wasting data efficiency.

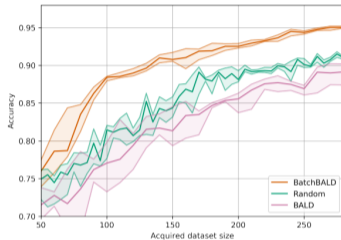


Figure 2: Performance on Repeated MNIST with acquisition size 10. See section 4.1 for further details. BatchBALD outperforms BALD while BALD performs worse than random acquisition due to the replications in the dataset.

(Credit: [Kirsch et al., 2019].)

- The problem with the “natural generalization”:

$$\text{acq}(f, B) = \sum_{\mathbf{x} \in B} \text{acq}(f, \mathbf{x}) \quad (78)$$

is that the sum doesn't consider the *overlap* between the information carried by different $\mathbf{x} \in b$.

- The problem with the “natural generalization”:

$$\text{acq}(f, B) = \sum_{\mathbf{x} \in B} \text{acq}(f, \mathbf{x}) \quad (78)$$

is that the sum doesn't consider the *overlap* between the information carried by different $\mathbf{x} \in B$.

Idea: don't break the acquisition function into a sum! For BALD, this means replacing:

$$\sum_{\mathbf{x} \in B} \left\{ \underbrace{H(Y | \mathbf{x}, L) - \mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta} | L)} [H(Y | \mathbf{x}, \boldsymbol{\theta})]}_{MI(Y, \Theta | \mathbf{x}, L)} \right\} \quad (79)$$

with

$$MI(\{Y_1, \dots, Y_b\}, \Theta | \{\mathbf{x}_1, \dots, \mathbf{x}_b\}, L) \quad (80)$$

- The problem with the “natural generalization”:

$$\text{acq}(f, B) = \sum_{\mathbf{x} \in B} \text{acq}(f, \mathbf{x}) \quad (78)$$

is that the sum doesn't consider the *overlap* between the information carried by different $\mathbf{x} \in B$.

Idea: don't break the acquisition function into a sum! For BALD, this means replacing:

$$\sum_{\mathbf{x} \in B} \left\{ \underbrace{H(Y | \mathbf{x}, L) - \mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta} | L)} [H(Y | \mathbf{x}, \boldsymbol{\theta})]}_{MI(Y, \Theta | \mathbf{x}, L)} \right\} \quad (79)$$

with

$$MI(\{Y_1, \dots, Y_b\}, \Theta | \{\mathbf{x}_1, \dots, \mathbf{x}_b\}, L) \quad (80)$$

- In other words, **don't assume independence** between the elements of B !

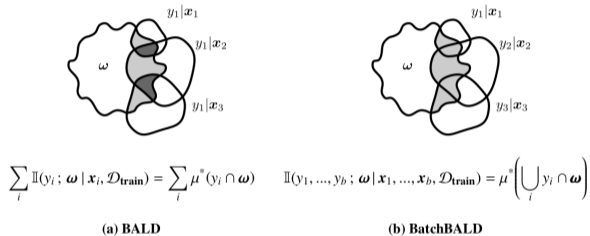


Figure 3: Intuition behind BALD and BatchBALD using *I*-diagrams [30]. *BALD* overestimates the joint mutual information. *BatchBALD*, however, takes the overlap between variables into account and will strive to acquire a better cover of ω . Areas contributing to the respective score are shown in grey, and areas that are double-counted in dark grey.

(Credit: [Kirsch et al., 2019].)

Problem: batch selection amounts to solving

$$\operatorname{argmax}_{B \subseteq U: |B|=b} MI(\{Y_1, \dots, Y_b\}, \Theta \mid \{x_1, \dots, x_b\}, L) \quad (81)$$

How can we solve this?

Problem: batch selection amounts to solving

$$\operatorname{argmax}_{B \subseteq U: |B|=b} MI(\{Y_1, \dots, Y_b\}, \Theta \mid \{x_1, \dots, x_b\}, L) \quad (81)$$

How can we solve this?

Submodular Function [Krause and Guestrin, 2008]

Let S be a set. A function f that maps subsets of S to real values is **submodular** if for every $B \subset A \subseteq S$ and any $x \in S \setminus A$ it holds that:

$$f(A \cup \{x\}) - f(A) \leq f(B \cup \{x\}) - f(B) \quad (82)$$

f enjoys a **diminishing returns** property: *adding an element x to a smaller set B “adds more” than adding the same element to a superset $A \supset B$.*

Problem: batch selection amounts to solving

$$\operatorname{argmax}_{B \subseteq U: |B|=b} MI(\{Y_1, \dots, Y_b\}, \Theta \mid \{x_1, \dots, x_b\}, L) \quad (81)$$

How can we solve this?

Submodular Function [Krause and Guestrin, 2008]

Let S be a set. A function f that maps subsets of S to real values is **submodular** if for every $B \subset A \subseteq S$ and any $x \in S \setminus A$ it holds that:

$$f(A \cup \{x\}) - f(A) \leq f(B \cup \{x\}) - f(B) \quad (82)$$

f enjoys a **diminishing returns** property: adding an element x to a smaller set B “adds more” than adding the same element to a superset $A \supset B$.

Maximizing a Submodular Functions

Let $f(A)$ be submodular and S the domain. Then, the greedy algorithm:

- $A_1 \leftarrow \emptyset$.
- $A_{t+1} \leftarrow \operatorname{argmax}_{x \in (S \setminus A_t)} f(A_t \cup \{x\})$.
- Stop when budget T is exhausted.

finds $A_T \subseteq S$ that has score $(1 - \frac{1}{e}) \approx 67\%$ as good as the score of the global optimum $A^* \subseteq S$ of f .

Problem: batch selection amounts to solving

$$\operatorname{argmax}_{B \subseteq U: |B|=b} MI(\{Y_1, \dots, Y_b\}, \Theta \mid \{\mathbf{x}_1, \dots, \mathbf{x}_b\}, L) \quad (83)$$

How can we solve this?

Problem: batch selection amounts to solving

$$\operatorname{argmax}_{B \subseteq U: |B|=b} MI(\{Y_1, \dots, Y_b\}, \Theta \mid \{x_1, \dots, x_b\}, L) \quad (83)$$

How can we solve this?

■ The mutual information is **submodular!** → apply greedy optimization:

- Pick x_1 to optimize $MI(Y_1, \Theta \mid x_1, L)$ (BALD)
- Pick x_{t+1} to optimize $MI(Y_{t+1} \cup Y_{1:t}, \Theta \mid x_{t+1} \cup x_{1:t}, L)$ (BALD over updated MI)

where $B = \{x_1, \dots, x_b\}$.

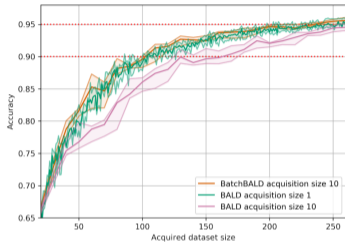


Figure 5: Performance on MNIST. BatchBALD outperforms BALD with acquisition size 10 and performs close to the optimum of acquisition size 1.

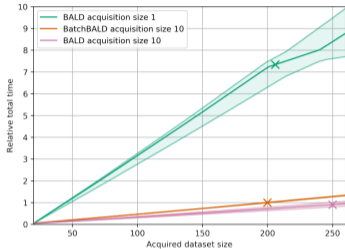


Figure 6: Relative total time on MNIST. Normalized to training BatchBALD with acquisition size 10 to 95% accuracy. The stars mark when 95% accuracy is reached for each method.

(Credit: [Kirsch et al., 2019].)

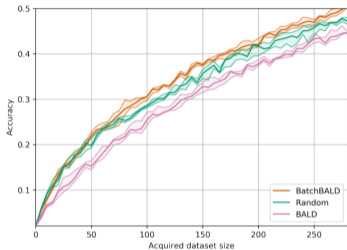


Figure 7: *Performance on EMNIST.* BatchBALD consistently outperforms both random acquisition and BALD while BALD is unable to beat random acquisition.

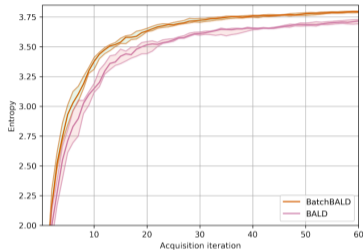


Figure 8: *Entropy of acquired class labels over acquisition steps on EMNIST.* BatchBALD steadily acquires a more diverse set of data points.

(Credit: [Kirsch et al., 2019].)

Idea: we want instances in L to be distributed as ground-truth

- This involves getting $p^*(\mathbf{X})$ right
- It also involves getting all the *modes* in it right

Idea: we want instances in L to be distributed as ground-truth

- This involves getting $p^*(\mathbf{X})$ right
- It also involves getting all the *modes* in it right

■ **Highly non-trivial:** we could train a generative model $\hat{p}_\theta(\mathbf{X}, Y)$ using **density estimation** and use that to guide query selection

Idea: we want instances in L to be distributed as ground-truth

- This involves getting $p^*(\mathbf{X})$ right
- It also involves getting all the *modes* in it right

■ **Highly non-trivial:** we could train a generative model $\hat{p}_\theta(\mathbf{X}, Y)$ using **density estimation** and use that to guide query selection → hard to train, break down in high dimension (in general)

Idea: we want instances in L to be distributed as ground-truth

- This involves getting $p^*(\mathbf{X})$ right
- It also involves getting all the *modes* in it right

■ **Highly non-trivial:** we could train a generative model $\hat{p}_\theta(\mathbf{X}, Y)$ using **density estimation** and use that to guide query selection → hard to train, break down in high dimension (in general)

■ Is there any alternative?

Idea: if we **cannot** find a classifier $f \in \mathcal{F}$ that tells L from U apart, and the latter is large enough (i.e., it can be used to approximate the ground-truth distribution $p^*(\mathbf{X})$), then L is high-quality.

Idea: if we **cannot** find a classifier $f \in \mathcal{F}$ that tells L from U apart, and the latter is large enough (i.e., it can be used to approximate the ground-truth distribution $p^*(\mathbf{X})$), then L is high-quality.

- Given $x \in U$, how certain are we that it comes from U rather than from L ?
 - if indistinguishable, we represented the true distribution using L
 - if distinguishable, it looks different from L so labeling it should be informative

\mathcal{F} -divergence

Given two distributions $p_S(X)$ and $p_T(X)$ on $X \in \mathcal{X}$ and a hypothesis class \mathcal{F} also on \mathcal{X} , the \mathcal{F} -divergence between p_S and p_T is:

$$d_{\mathcal{F}}(p_S, p_T) = 2 \cdot \sup_{f \in \mathcal{F}} |p_S(\{x : f(x) = 1\}) - p_T(\{x : f(x) = 1\})| \quad (84)$$

\mathcal{F} -divergence

Given two distributions $p_S(X)$ and $p_T(X)$ on $X \in \mathcal{X}$ and a hypothesis class \mathcal{F} also on \mathcal{X} , the \mathcal{F} -divergence between p_S and p_T is:

$$d_{\mathcal{F}}(p_S, p_T) = 2 \cdot \sup_{f \in \mathcal{F}} |p_S(\{x : f(x) = 1\}) - p_T(\{x : f(x) = 1\})| \quad (84)$$

- Measures how different two domains p_S and p_T are from the perspective of a model class \mathcal{F} : the larger the difference, the more different they look.
- \mathcal{F} -divergence used to identify **concept drift** and – symmetrically – to estimate how well a classifier trained on one task will perform on a different, related task

Idea: use \mathcal{F} -divergence to discriminate between L and U

Idea: use \mathcal{F} -divergence to discriminate between L and U

■ The \mathcal{F} -divergence is:

$$d_{\mathcal{F}}(p_S, p_T) = 2 \cdot \sup_{f \in \mathcal{F}} |p_S(\{x : f(x) = 1\}) - p_T(\{x : f(x) = 1\})| \quad (85)$$

Idea: use \mathcal{F} -divergence to discriminate between L and U

■ The \mathcal{F} -divergence is:

$$d_{\mathcal{F}}(p_S, p_T) = 2 \cdot \sup_{f \in \mathcal{F}} |p_S(\{x : f(x) = 1\}) - p_T(\{x : f(x) = 1\})| \quad (85)$$

■ Approximate as follows:

- $p_S(\mathbf{x}) := \frac{1}{|L|} \sum_{\mathbf{x}' \in L} \delta\{\mathbf{x}' = \mathbf{x}\}$
- $p_T(\mathbf{x}) := \frac{1}{|U|} \sum_{\mathbf{x}' \in U} \delta\{\mathbf{x}' = \mathbf{x}\}$.

Then \mathcal{F} -divergence becomes:

$$2 \sup_{f \in \mathcal{F}} \left| \frac{1}{|L|} \sum_{\mathbf{x} \in L} f(\mathbf{x}) - \frac{1}{|U|} \sum_{\mathbf{x} \in U} f(\mathbf{x}) \right| \quad (86)$$

Idea: use \mathcal{F} -divergence to discriminate between L and U

■ The \mathcal{F} -divergence is:

$$d_{\mathcal{F}}(p_S, p_T) = 2 \cdot \sup_{f \in \mathcal{F}} |p_S(\{x : f(x) = 1\}) - p_T(\{x : f(x) = 1\})| \quad (85)$$

■ Approximate as follows:

- $p_S(\mathbf{x}) := \frac{1}{|L|} \sum_{\mathbf{x}' \in L} \delta\{\mathbf{x}' = \mathbf{x}\}$
- $p_T(\mathbf{x}) := \frac{1}{|U|} \sum_{\mathbf{x}' \in U} \delta\{\mathbf{x}' = \mathbf{x}\}$.

Then \mathcal{F} -divergence becomes:

$$2 \sup_{f \in \mathcal{F}} \left| \frac{1}{|L|} \sum_{\mathbf{x} \in L} f(\mathbf{x}) - \frac{1}{|U|} \sum_{\mathbf{x} \in U} f(\mathbf{x}) \right| \quad (86)$$

■ How to compute this?

Idea: use \mathcal{F} -divergence to discriminate between L and U

■ The \mathcal{F} -divergence is:

$$d_{\mathcal{F}}(p_S, p_T) = 2 \cdot \sup_{f \in \mathcal{F}} |p_S(\{x : f(x) = 1\}) - p_T(\{x : f(x) = 1\})| \quad (85)$$

■ Approximate as follows:

- $p_S(\mathbf{x}) := \frac{1}{|L|} \sum_{\mathbf{x}' \in L} \delta\{\mathbf{x}' = \mathbf{x}\}$
- $p_T(\mathbf{x}) := \frac{1}{|U|} \sum_{\mathbf{x}' \in U} \delta\{\mathbf{x}' = \mathbf{x}\}$.

Then \mathcal{F} -divergence becomes:

$$2 \sup_{f \in \mathcal{F}} \left| \frac{1}{|L|} \sum_{\mathbf{x} \in L} f(\mathbf{x}) - \frac{1}{|U|} \sum_{\mathbf{x} \in U} f(\mathbf{x}) \right| \quad (86)$$

■ How to compute this? $\sup_{f \in \mathcal{F}}$ can be implemented by learning f from data set:

- Given L and U , define binary classification task:

- Given L and U , define binary classification task:
 - For all $x \in L$, add $(\phi(x), \text{labeled})$ to dataset

- Given L and U , define binary classification task:
 - For all $\mathbf{x} \in L$, add $(\phi(\mathbf{x}), \text{labeled})$ to dataset
 - For all $\mathbf{x} \in U$, add $(\phi(\mathbf{x}), \text{unlabeled})$ to dataset

- Given L and U , define binary classification task:
 - For all $\mathbf{x} \in L$, add $(\phi(\mathbf{x}), \text{labeled})$ to dataset
 - For all $\mathbf{x} \in U$, add $(\phi(\mathbf{x}), \text{unlabeled})$ to dataset
 - Train (simple) classifier $p_\psi(\cdot \mid \phi(\mathbf{x}))$ to distinguish between the two sources by optimizing **cross-entropy loss**

- Given L and U , define binary classification task:
 - For all $\mathbf{x} \in L$, add $(\phi(\mathbf{x}), \text{labeled})$ to dataset
 - For all $\mathbf{x} \in U$, add $(\phi(\mathbf{x}), \text{unlabeled})$ to dataset
 - Train (simple) classifier $p_\psi(\cdot \mid \phi(\mathbf{x}))$ to distinguish between the two sources by optimizing **cross-entropy loss**

This implicitly makes use of the labels y in L through $\phi(\mathbf{x})$.

- Given L and U , define binary classification task:
 - For all $\mathbf{x} \in L$, add $(\phi(\mathbf{x}), \text{labeled})$ to dataset
 - For all $\mathbf{x} \in U$, add $(\phi(\mathbf{x}), \text{unlabeled})$ to dataset
 - Train (simple) classifier $p_\psi(\cdot \mid \phi(\mathbf{x}))$ to distinguish between the two sources by optimizing **cross-entropy loss**

This implicitly makes use of the labels y in L through $\phi(\mathbf{x})$.

- Pick those instances $\mathbf{x} \in U$ that have lowest probability $p_\psi(\text{labeled} \mid \phi(\mathbf{x}))$

- Given L and U , define binary classification task:
 - For all $\mathbf{x} \in L$, add $(\phi(\mathbf{x}), \text{labeled})$ to dataset
 - For all $\mathbf{x} \in U$, add $(\phi(\mathbf{x}), \text{unlabeled})$ to dataset
 - Train (simple) classifier $p_\psi(\cdot \mid \phi(\mathbf{x}))$ to distinguish between the two sources by optimizing **cross-entropy loss**

This implicitly makes use of the labels y in L through $\phi(\mathbf{x})$.

- Pick those instances $\mathbf{x} \in U$ that have lowest probability $p_\psi(\text{labeled} \mid \phi(\mathbf{x}))$

Remark: this is closely related to GANs!

Conclusion and Further Reading

- AL useful when supervision is expensive high → choose it wisely
- Many variants: pool-based, streaming, and query synthesis
- Many practical approaches: **uncertainty**-based (uncertainty sampling, QBC, expected gradient length), **diversity**-based (information density).

Some can be derived from *version spaces* and *model improvement*.

- Deep variants select entire **batches** and often rely on Bayesian techniques
- Critique & realistic annotators, costs, etc.: [Herde et al., 2021] [Settles, 2011]
- Plenty of room for new research ;-)



Baum, E. B. and Lang, K. (1992).

Query learning can work poorly when a human oracle is used.

In *International joint conference on neural networks*, volume 8, page 8.



Gal, Y. and Ghahramani, Z. (2016).

Dropout as a bayesian approximation: Representing model uncertainty in deep learning.

In *international conference on machine learning*, pages 1050–1059. PMLR.



Gal, Y., Islam, R., and Ghahramani, Z. (2017).

Deep bayesian active learning with image data.

In *International Conference on Machine Learning*, pages 1183–1192. PMLR.



Herde, M., Huseljic, D., Sick, B., and Calma, A. (2021).

A survey on cost types, interaction schemes, and annotator performance models in selection algorithms for active learning in classification.

arXiv preprint arXiv:2109.11301.



Hüllermeier, E. and Waegeman, W. (2021).

Aleatoric and epistemic uncertainty in machine learning: An introduction to concepts and methods.

Machine Learning, 110(3):457–506.



King, R. D., Rowland, J., Oliver, S. G., Young, M., Aubrey, W., Byrne, E., Liakata, M., Markham, M., Pir, P., Soldatova, L. N., et al. (2009).

The automation of science.

Science, 324(5923):85–89.



Kirsch, A., Van Amersfoort, J., and Gal, Y. (2019).

Batchbald: Efficient and diverse batch acquisition for deep bayesian active learning.

Advances in neural information processing systems, 32:7026–7037.



Krause, A. and Guestrin, C. (2008).

Beyond convexity: Submodularity in machine learning.

ICML Tutorials.



Kristiadi, A., Hein, M., and Hennig, P. (2020).

Being bayesian, even just a bit, fixes overconfidence in relu networks.

In *International Conference on Machine Learning*, pages 5436–5446. PMLR.



Nguyen, A., Dosovitskiy, A., Yosinski, J., Brox, T., and Clune, J. (2016).

Synthesizing the preferred inputs for neurons in neural networks via deep generator networks.

Advances in neural information processing systems, 29:3387–3395.



Settles, B. (2011).

From theories to queries: Active learning in practice.

In *Active Learning and Experimental Design workshop In conjunction with AISTATS 2010*, pages 1–18.

JMLR Workshop and Conference Proceedings.



Settles, B. (2012).

Active learning.



Shalev-Shwartz, S. and Ben-David, S. (2014).

Understanding machine learning: From theory to algorithms.

Cambridge university press.



Sharma, M. and Bilgic, M. (2017).

Evidence-based uncertainty sampling for active learning.

Data Mining and Knowledge Discovery, 31(1):164–202.