

EMPIRICAL RISK MINIMIZATION

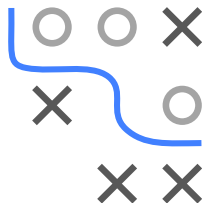
Very often, in ML, we minimize the empirical risk

$$\mathcal{R}_{\text{emp}}(f) = \sum_{i=1}^n L(y^{(i)}, f(\mathbf{x}^{(i)}))$$

- each observation $(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{X} \times \mathcal{Y}$, so from feature and target space
- $f_{\mathcal{H}} : \mathcal{X} \rightarrow \mathbb{R}^g$, f is a model from hypothesis space \mathcal{H} ; maps a feature vector to output score; sometimes or often we omit \mathcal{H} in the index
- $L : (\mathcal{Y} \times \mathbb{R}^g) \rightarrow \mathbb{R}$ is loss;
 $L(y, f)$ measures distance between label and prediction
- We assume that $(\mathbf{x}, y) \sim \mathbb{P}_{xy}$ and $(\mathbf{x}^{(i)}, y^{(i)}) \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_{xy}$
 \mathbb{P}_{xy} is the distribution of the data generating process (DGP)

Let's define (and minimize) loss in expectation, the theoretical risk:

$$\mathcal{R}(f) := \mathbb{E}_{xy}[L(y, f(\mathbf{x}))] = \int L(y, f(\mathbf{x})) d\mathbb{P}_{xy}$$

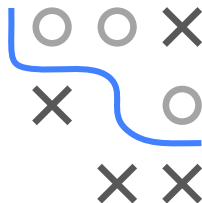


TWO SHORT EXAMPLES

Regression with linear model:

- Model: $f(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{x} + \theta_0$
- Squared loss: $L(y, f) = (y - f)^2$
- Hypothesis space:

$$\mathcal{H}_{\text{lin}} = \left\{ \mathbf{x} \mapsto \boldsymbol{\theta}^\top \mathbf{x} + \theta_0 : \boldsymbol{\theta} \in \mathbb{R}^d, \theta_0 \in \mathbb{R} \right\}$$



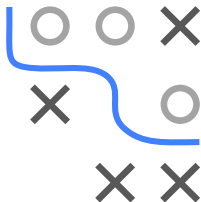
Binary classification with shallow MLP:

- Model: $f(\mathbf{x}) = \pi(\mathbf{x}) = \sigma(\mathbf{w}_2^\top \text{ReLU}(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + b_2)$
- Binary cross-entropy loss:
 $L(y, \pi) = -(y \log(\pi) + (1 - y) \log(1 - \pi))$
- Hypothesis space:

$$\mathcal{H}_{\text{MLP}} = \left\{ \mathbf{x} \mapsto \sigma(\mathbf{w}_2^\top \text{ReLU}(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + b_2) : \mathbf{W}_1 \in \mathbb{R}^{h \times d}, \mathbf{b}_1 \in \mathbb{R}^h, \mathbf{w}_2 \in \mathbb{R}^h, b_2 \in \mathbb{R} \right\}$$

OPTIMAL CONSTANTS FOR A LOSS

- Let's assume some RV $z \in \mathcal{Y}$ for a label
- z not RV y , because we want to fiddle with its distribution
- Assume z has distribution Q , so $z \sim Q$
- We can now consider $\arg \min_c \mathbb{E}_{z \sim Q}[L(z, c)]$
so the score-constant which loss-minimally approximates z



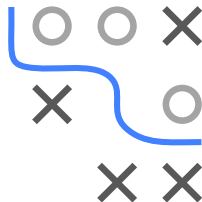
We will consider 3 cases for Q

- $Q = P_y$, simply our labels and their marginal distribution in \mathbb{P}_{xy}
- $Q = P_{y|x=\tilde{x}}$, conditional label distribution at point $x = \tilde{x}$
- $Q = P_n$, the empirical product distribution for data y_1, \dots, y_n

If we can solve $\arg \min_c \mathbb{E}_{z \sim Q}[L(z, c)]$ for any Q , we will get multiple useful results!

OPTIMAL CONSTANT MODEL

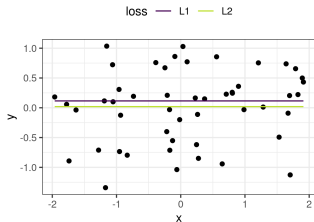
- We would like a loss optimal, constant baseline predictor
- A "featureless" ML model, which always predicts the same value
- Can use it as baseline in experiments, if we don't beat this with more complex model, that model is useless
- Will also be useful as component in algorithms and derivations



$$f_c^* = \arg \min_{c \in \mathbb{R}} \mathbb{E}_{xy} [L(y, c)] = \arg \min_{c \in \mathbb{R}} \mathbb{E}_y [L(y, c)]$$

and $f(\mathbf{x}) = \theta = c$ that optimizes the empirical risk $\mathcal{R}_{\text{emp}}(\theta)$ is denoted as as

$$\hat{f}_c = \arg \min_{c \in \mathbb{R}} \mathcal{R}_{\text{emp}}(\theta).$$



OPTIMAL CONSTANT MODEL

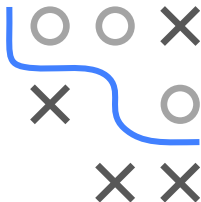
- Let's start with the simplest case, L2 loss
- And we want to find an optimal constant model for

$$\begin{aligned}\arg \min \mathbb{E}[L(z, c)] &= \\ \arg \min \mathbb{E}[(z - c)^2] &= \\ \arg \min \mathbb{E}[z^2] - 2c\mathbb{E}[z] + c^2 &= \\ E[z] &\end{aligned}$$

- Using $Q = P_y$, this means that, given we know the label distribution, the best constant is $c = E[y]$.
- If we only have data y_1, \dots, y_n

$$\arg \min \mathbb{E}_{z \sim P_n}[(z - c)^2] = \mathbb{E}_{z \sim P_n}[z] = \frac{1}{n} \sum_{i=1}^n y^{(i)} = \bar{y}$$

- And we want to find an optimal constant model for

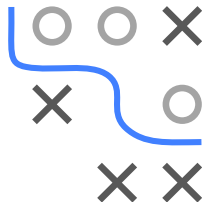


RISK MINIMIZER

Let us assume we are in an “ideal world”:

- The hypothesis space $\mathcal{H} = \mathcal{H}_{all}$ is unrestricted. We can choose any measurable $f : \mathcal{X} \rightarrow \mathbb{R}^g$.
- We also assume an ideal optimizer; the risk minimization can always be solved perfectly and efficiently.
- We know \mathbb{P}_{xy} .

How should f be chosen?



RISK MINIMIZER / 2

The f with minimal risk across all (measurable) functions is called the **risk minimizer**, **population minimizer** or **Bayes optimal model**.

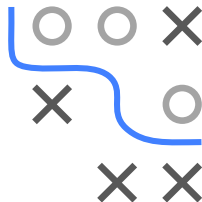
$$\begin{aligned} f_{\mathcal{H}_{all}}^* &= \arg \min_{f \in \mathcal{H}_{all}} \mathcal{R}(f) = \arg \min_{f \in \mathcal{H}_{all}} \mathbb{E}_{xy} [L(y, f(\mathbf{x}))] \\ &= \arg \min_{f \in \mathcal{H}_{all}} \int L(y, f(\mathbf{x})) d\mathbb{P}_{xy}. \end{aligned}$$

The resulting risk is called **Bayes risk**: $\mathcal{R}^* = \mathcal{R}(f_{\mathcal{H}_{all}}^*)$

Note that if we leave out the hypothesis space in the subscript it becomes clear from the context!

Similarly, we define the risk minimizer over some $\mathcal{H} \subset \mathcal{H}_{all}$ as

$$f_{\mathcal{H}}^* = \arg \min_{f \in \mathcal{H}} \mathcal{R}(f)$$



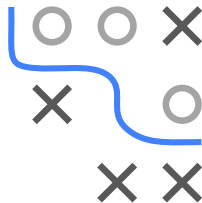
OPTIMAL POINT-WISE PREDICTIONS

To derive the risk minimizer, observe that by law of total expectation

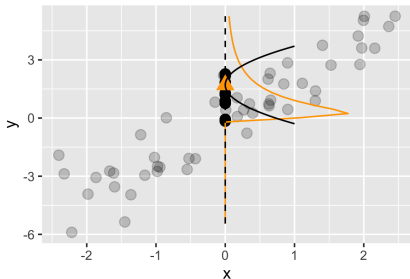
$$\mathcal{R}(f) = \mathbb{E}_{xy} [L(y, f(\mathbf{x}))] = \mathbb{E}_x [\mathbb{E}_{y|x} [L(y, f(\mathbf{x})) \mid \mathbf{x}]] .$$

- We can choose $f(\mathbf{x})$ as we want (unrestricted hypothesis space, no assumed functional form)
- Hence, for a fixed value $\mathbf{x} \in \mathcal{X}$ we can select **any** value c we want to predict. So we construct the **point-wise optimizer**

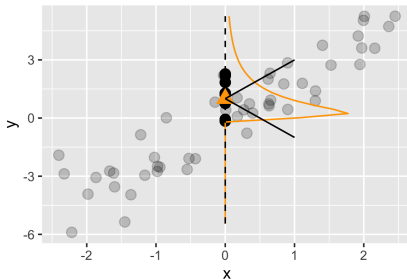
$$f^*(\tilde{\mathbf{x}}) = \operatorname{argmin}_c \mathbb{E}_{y|x} [L(y, c) \mid \mathbf{x} = \tilde{\mathbf{x}}]$$



L2 Loss: Fix one x



L1 Loss: Fix one x



THEORETICAL AND EMPIRICAL RISK

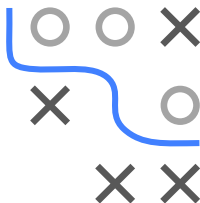
The risk minimizer is mainly a theoretical tool:

- In practice we need to restrict the hypothesis space \mathcal{H} such that we can efficiently search over it.
- In practice we (usually) do not know \mathbb{P}_{xy} . Instead of $\mathcal{R}(f)$, we are optimizing the empirical risk

$$\hat{f}_{\mathcal{H}} = \arg \min_{f \in \mathcal{H}} \mathcal{R}_{\text{emp}}(f) = \arg \min_{f \in \mathcal{H}} \sum_{i=1}^n L(y^{(i)}, f(\mathbf{x}^{(i)}))$$

Note that according to the **law of large numbers** (LLN), the empirical risk converges to the true risk (but beware of overfitting!):

$$\bar{\mathcal{R}}_{\text{emp}}(f) = \frac{1}{n} \sum_{i=1}^n L(y^{(i)}, f(\mathbf{x}^{(i)})) \xrightarrow{n \rightarrow \infty} \mathcal{R}(f).$$



ESTIMATION AND APPROXIMATION ERROR

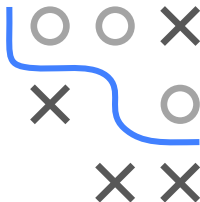
Goal of learning: Train a model $\hat{f}_{\mathcal{H}}$ for which the true risk $\mathcal{R}(\hat{f}_{\mathcal{H}})$ is close to the Bayes risk \mathcal{R}^* . In other words, we want the **Bayes regret** or **excess risk**

$$\mathcal{R}(\hat{f}_{\mathcal{H}}) - \mathcal{R}^*$$

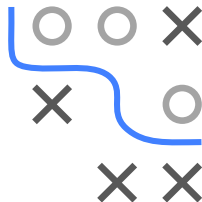
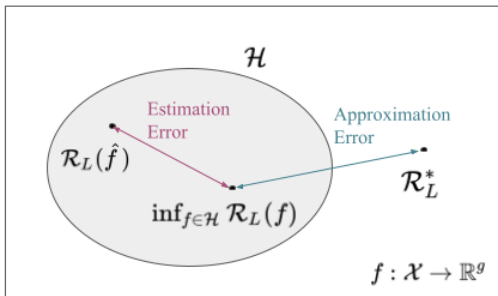
to be as low as possible.

The Bayes regret can be decomposed as follows:

$$\begin{aligned}\mathcal{R}(\hat{f}_{\mathcal{H}}) - \mathcal{R}^* &= \underbrace{\left[\mathcal{R}(\hat{f}_{\mathcal{H}}) - \inf_{f \in \mathcal{H}} \mathcal{R}(f) \right]}_{\text{estimation error}} + \underbrace{\left[\inf_{f \in \mathcal{H}} \mathcal{R}(f) - \mathcal{R}^* \right]}_{\text{approximation error}} \\ &= \left[\mathcal{R}(\hat{f}_{\mathcal{H}}) - \mathcal{R}(f_{\mathcal{H}}^*) \right] + \left[\mathcal{R}(f_{\mathcal{H}}^*) - \mathcal{R}(f_{\mathcal{H}_{all}}^*) \right]\end{aligned}$$



ESTIMATION AND APPROXIMATION ERROR / 2



- $\mathcal{R}(\hat{f}) - \inf_{f \in \mathcal{H}} \mathcal{R}(f)$ is the **estimation error**. We fit \hat{f} via empirical risk minimization and (usually) use approximate optimization, so we usually do not find the optimal $f \in \mathcal{H}$.
- $\inf_{f \in \mathcal{H}} \mathcal{R}(f) - \mathcal{R}^*$ is the **approximation error**. We need to restrict to a hypothesis space \mathcal{H} which might not even contain the Bayes optimal model f^* .

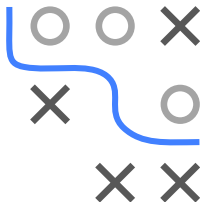
(UNIVERSALLY) CONSISTENT LEARNERS

Consistency is an asymptotic property of a learning algorithm, which ensures the algorithm returns **the correct model** when given **unlimited data**.

Let $\mathcal{I} : \mathbb{D} \rightarrow \mathcal{H}$ be a learning algorithm that takes a training set $\mathcal{D}_{\text{train}} \sim \mathbb{P}_{xy}$ of size n_{train} and estimates a model $\hat{f} : \mathcal{X} \rightarrow \mathbb{R}^g$.

The learning method \mathcal{I} is said to be **consistent** w.r.t. a certain distribution \mathbb{P}_{xy} if the risk of the estimated model \hat{f} converges in probability (“ \xrightarrow{p} ”) to the Bayes risk \mathcal{R}^* when n_{train} goes to ∞ :

$$\mathcal{R}(\mathcal{I}(\mathcal{D}_{\text{train}})) \xrightarrow{p} \mathcal{R}^* \quad \text{for } n_{\text{train}} \rightarrow \infty.$$



(UNIVERSALLY) CONSISTENT LEARNERS / 2

Consistency is defined w.r.t. a particular distribution \mathbb{P}_{xy} . But since we usually do not know \mathbb{P}_{xy} , consistency does not offer much help to choose an algorithm for a particular task.

More interesting is the stronger concept of **universal consistency**: An algorithm is universally consistent if it is consistent for **any** distribution.

In Stone's famous consistency theorem from 1977, the universal consistency of a weighted average estimator as KNN was proven. Many other ML models have since then been proven to be universally consistent (SVMs, ANNs, etc.).

Note that universal consistency is obviously a desirable property - however, (universal) consistency does not tell us anything about convergence rates ...

