# Introduction to Machine Learning
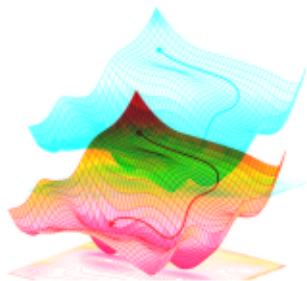
## Advanced Risk Minimization

## Risk Minimizers



### Learning goals

- Bayes optimal model (also: risk minimizer, population minimizer)
- Bayes risk
- Bayes regret, estimation and approximation error
- Optimal constant model
- Consistent learners

# EMPIRICAL RISK MINIMIZATION

Very often, in ML, we minimize the empirical risk

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

where

- each observation $\left(\mathbf{x}^{(i)}, y^{(i)}\right) \in \mathcal{X} \times \mathcal{Y}$, so from feature and target space
- $h_{\mathcal{L}} : \mathcal{X} \to \mathbb{R}^g$, $f$ 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
- $L : (\mathcal{Y} \times \mathbb{R}^g) \to \mathbb{R}$ is loss; $L(y, f)$ measures distance between label and prediction
- We assume that $(\mathbf{x}, y) \sim P_{xy}$ and $(\mathbf{x}^{(i)}, y^{(i)}) \overset{\text{i.i.d.}}{\sim} P_{xy}$; $P_{xy}$ is the distribution of the data generating process (DGP)

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

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

What is the theoretical justification for this procedure?

# TWO SHORT EXAMPLES

**Regression with linear model:**

- Model: $f(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{x} + \theta_0$
- Squared loss: $L(y, f) = (y - f(\mathbf{x}))^2$
- 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 $P_{xy}$
- $Q = P_{y|x=x}$, conditional label distribution at point $x = \tilde{x}$
- $Q = P_n$, the empirical product distribution for data $y_1, \ldots, y_n$

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

- 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
- Hence, for a fixed value $x \in \mathcal{X}$ we can select **any** value $c$ we want to predict

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

and $f(x) = \theta = c$ that optimizes the empirical risk $\mathcal{R}_{emp}(\theta)$ is denoted as as
$\hat{f}_c = \arg\min_{c\in\mathbb{R}} \mathcal{R}_{emp}(\theta)$.

L1 Loss: Fix one x

# OPTIMAL CONSTANT MODEL

- Let's start with the simplest case, L2 loss
- And we want to find and optimal constant model for $\mathcal{H}$ such that we can efficiently search over it.
- In practice we (usually) do not know. Instead of $\mathcal{R}(f)$, we are optimizing the empirical risk

$$\arg\min E[L(z,c)] =$$

$$\arg\min E[(z-c)^2] =$$

$$\arg\min E[z^2] - 2cE[z] + c^2 =$$

$$E[z]$$

- 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 $y_1, \ldots y_n$

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

- And we want to find and optimal constant model for

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.

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

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

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

# ...D APPROXIMATION ERROR

**Goal of learning:** Train a model $\hat{f}$ for which the true risk $\mathcal{R}_L\left(\hat{f}\right)$ is

Let us assume we are in an "ideal world":

- The hypothesis space $\mathcal{H} = \mathcal{H}_{all}$ is unrestricted. We can choose any measurable $f : \mathcal{X} \to \mathbb{R}^g$.

- We also assume an ideal optimizer; the risk minimization can always be solved perfectly and efficiently.

- We know $\mathcal{R}_L$ can be decomposed as follows:

How should $f$ be chosen?

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

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

$$
\begin{aligned}
f^*_{\mathcal{H}_{all}} &= \arg\min_{f \in \mathcal{H}_{all}} \mathcal{R}(f) = \arg\min_{f \in \mathcal{H}_{all}} \mathbb{E}_{xy}[L(y, f(\mathbf{x}))] \\
&= \arg\min_{f \in \mathcal{H}_{all}} \int L(y, f(\mathbf{x})) \, dP_{xy} \\
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad f : \mathcal{X} \to \mathbb{R}^g
\end{aligned}
$$

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

- $\mathcal{R}_L\left(\hat{f}\right) - \inf_{f \in \mathcal{H}} \mathcal{R}(f)$ 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 $f \in \mathcal{H}$.

- $\inf_{f \in \mathcal{H}} \mathcal{R}_L(f) - \mathcal{R}_L^*$ is the **approximation error**. We need to restrict to a hypothesis space $\mathcal{H}$ which might not even contain the Bayes optimal model.

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

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

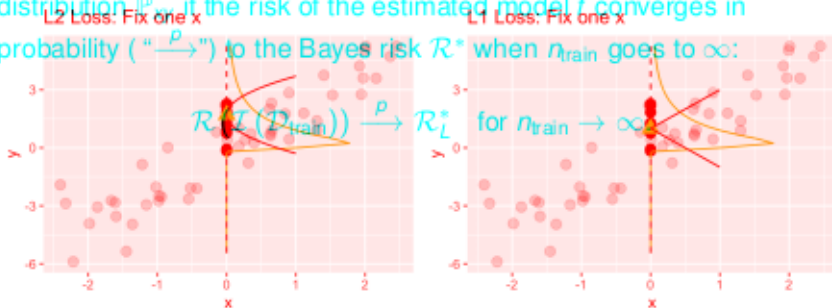$$\mathcal{R}(f) = \mathbb{E}_{xy}[L(y, f(\mathbf{x}))] = \mathbb{E}_x[\mathbb{E}_{y|x}[L(y, f(\mathbf{x})) \mid \mathbf{x}]] \ .$$

- 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**

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

L2 Loss: Fix one x      L1 Loss: Fix one x

# THEORETICAL AND EMPIRICAL RISKS / 2

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 $P_{xy}$. Instead of $\mathcal{R}(f)$, we are optimizing the empirical risk

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

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

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

# ESTIMATION AND APPROXIMATION ERROR

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

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

to be as low as possible.

The Bayes regret can be decomposed as follows:

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

$$\stackrel{.}{=} \left[\mathcal{R}(\hat{f}_\mathcal{H}) - \mathcal{R}(f^*_\mathcal{H})\right] + \left[\mathcal{R}(f^*_\mathcal{H}) - \mathcal{R}(f^*_{\mathcal{H}_{all}})\right]$$