

Deep learning

2.1. Loss and risk

François Fleuret

<https://fleuret.org/dlc/>

Dec 20, 2020

The general objective of machine learning is to capture regularity in data to make predictions.

In our regression example, we modeled age and blood pressure as being linearly related, to predict the latter from the former.

The general objective of machine learning is to capture regularity in data to make predictions.

In our regression example, we modeled age and blood pressure as being linearly related, to predict the latter from the former.

There are multiple types of inference that we can roughly split into three categories:

- Classification (e.g. object recognition, cancer detection, speech processing),
- regression (e.g. customer satisfaction, stock prediction, epidemiology), and
- density estimation (e.g. outlier detection, data visualization, sampling/synthesis).

The standard formalization considers a measure of probability

$$\mu_{X,Y}$$

over the observation/value of interest, and i.i.d. training samples

$$(x_n, y_n), \quad n = 1, \dots, N.$$

Intuitively, for classification it can often be interpreted as

$$\mu_{X,Y}(x,y) = \mu_{X|Y=y}(x) P(Y=y)$$

that is, draw Y first, and given its value, generate X .

Intuitively, for classification it can often be interpreted as

$$\mu_{X,Y}(x,y) = \mu_{X|Y=y}(x) P(Y=y)$$

that is, draw Y first, and given its value, generate X .

So the quantity

$$\mu_{X|Y=y}$$

stands for the distribution of the observable signal for the class y (e.g. “sound of an /ē/”, “image of a cat”).

For regression, one would interpret the joint law more naturally as

$$\mu_{X,Y}(x,y) = \mu_{Y|X=x}(y) \mu_X(x)$$

which would be: first, generate X , and given its value, generate Y .

For regression, one would interpret the joint law more naturally as

$$\mu_{X,Y}(x,y) = \mu_{Y|X=x}(y) \mu_X(x)$$

which would be: first, generate X , and given its value, generate Y .

In the simple cases

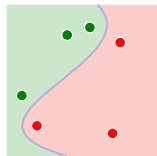
$$Y = f(X) + \epsilon$$

where f is the deterministic dependency between x and y , and ϵ is a random noise, independent of X .

With such a probabilistic perspective, we can more precisely define the three types of inferences we introduced before:

Classification,

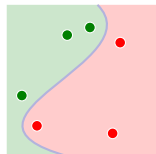
- (X, Y) random variables on $\mathcal{X} = \mathbb{R}^D \times \{1, \dots, C\}$,
- we want to estimate $\operatorname{argmax}_y P(Y = y \mid X = x)$.



With such a probabilistic perspective, we can more precisely define the three types of inferences we introduced before:

Classification,

- (X, Y) random variables on $\mathcal{X} = \mathbb{R}^D \times \{1, \dots, C\}$,
- we want to estimate $\operatorname{argmax}_y P(Y = y \mid X = x)$.



Regression,

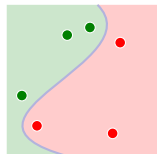
- (X, Y) random variables on $\mathcal{X} = \mathbb{R}^D \times \mathbb{R}$,
- we want to estimate $\mathbb{E}(Y \mid X = x)$.



With such a probabilistic perspective, we can more precisely define the three types of inferences we introduced before:

Classification,

- (X, Y) random variables on $\mathcal{X} = \mathbb{R}^D \times \{1, \dots, C\}$,
- we want to estimate $\operatorname{argmax}_y P(Y = y \mid X = x)$.



Regression,

- (X, Y) random variables on $\mathcal{X} = \mathbb{R}^D \times \mathbb{R}$,
- we want to estimate $\mathbb{E}(Y \mid X = x)$.



Density estimation,

- X random variable on $\mathcal{X} = \mathbb{R}^D$,
- we want to estimate μ_X .



The boundaries between these categories are fuzzy:

- Regression allows to do classification through class scores.
- Density models allow to do classification thanks to Bayes' law.

etc.

Risk, empirical risk

Learning consists of finding in a set \mathcal{F} of functionals a “good” f^* (or its parameters’ values) usually defined through a loss

$$\ell : \mathcal{F} \times \mathcal{Z} \rightarrow \mathbb{R}$$

such that $\ell(f, z)$ increases with how wrong f is on z .

Learning consists of finding in a set \mathcal{F} of functionals a “good” f^* (or its parameters’ values) usually defined through a loss

$$\ell : \mathcal{F} \times \mathcal{Z} \rightarrow \mathbb{R}$$

such that $\ell(f, z)$ increases with how wrong f is on z . For instance

- for classification:

$$\ell(f, (x, y)) = \mathbf{1}_{\{f(x) \neq y\}},$$

- for regression:

$$\ell(f, (x, y)) = (f(x) - y)^2,$$

- for density estimation:

$$\ell(q, z) = -\log q(z).$$

Learning consists of finding in a set \mathcal{F} of functionals a “good” f^* (or its parameters’ values) usually defined through a loss

$$\ell : \mathcal{F} \times \mathcal{Z} \rightarrow \mathbb{R}$$

such that $\ell(f, z)$ increases with how wrong f is on z . For instance

- for classification:

$$\ell(f, (x, y)) = \mathbf{1}_{\{f(x) \neq y\}},$$

- for regression:

$$\ell(f, (x, y)) = (f(x) - y)^2,$$

- for density estimation:

$$\ell(q, z) = -\log q(z).$$

The loss may include additional terms related to f itself.

We are looking for an f with a small **expected risk**

$$R(f) = \mathbb{E}_Z (\ell(f, Z)),$$

which means that our learning procedure would ideally choose

$$f^* = \underset{f \in \mathcal{F}}{\operatorname{argmin}} R(f).$$

We are looking for an f with a small **expected risk**

$$R(f) = \mathbb{E}_Z (\ell(f, Z)),$$

which means that our learning procedure would ideally choose

$$f^* = \underset{f \in \mathcal{F}}{\operatorname{argmin}} R(f).$$

Although this quantity is unknown, if we have i.i.d. training samples

$$\mathcal{D} = \{Z_1, \dots, Z_N\},$$

we can compute an estimate, the **empirical risk**:

$$\hat{R}(f; \mathcal{D}) = \hat{\mathbb{E}}_{\mathcal{D}}(\ell(f, Z)) = \frac{1}{N} \sum_{n=1}^N \ell(f, Z_n).$$

We have

$$\mathbb{E}_{Z_1, \dots, Z_N} \left(\hat{R}(f; \mathcal{D}) \right) = \mathbb{E}_{Z_1, \dots, Z_N} \left(\frac{1}{N} \sum_{n=1}^N \ell(f, Z_n) \right)$$

We have

$$\begin{aligned}\mathbb{E}_{Z_1, \dots, Z_N} \left(\hat{R}(f; \mathcal{D}) \right) &= \mathbb{E}_{Z_1, \dots, Z_N} \left(\frac{1}{N} \sum_{n=1}^N \ell(f, Z_n) \right) \\ &= \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{Z_n} (\ell(f, Z_n))\end{aligned}$$

We have

$$\begin{aligned}\mathbb{E}_{Z_1, \dots, Z_N} \left(\hat{R}(f; \mathcal{D}) \right) &= \mathbb{E}_{Z_1, \dots, Z_N} \left(\frac{1}{N} \sum_{n=1}^N \ell(f, Z_n) \right) \\ &= \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{Z_n} (\ell(f, Z_n)) \\ &= \mathbb{E}_Z (\ell(f, Z))\end{aligned}$$

We have

$$\begin{aligned}\mathbb{E}_{Z_1, \dots, Z_N} \left(\hat{R}(f; \mathcal{D}) \right) &= \mathbb{E}_{Z_1, \dots, Z_N} \left(\frac{1}{N} \sum_{n=1}^N \ell(f, Z_n) \right) \\ &= \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{Z_n} (\ell(f, Z_n)) \\ &= \mathbb{E}_Z (\ell(f, Z)) \\ &= R(f).\end{aligned}$$

We have

$$\begin{aligned}\mathbb{E}_{Z_1, \dots, Z_N} \left(\hat{R}(f; \mathcal{D}) \right) &= \mathbb{E}_{Z_1, \dots, Z_N} \left(\frac{1}{N} \sum_{n=1}^N \ell(f, Z_n) \right) \\ &= \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{Z_n} (\ell(f, Z_n)) \\ &= \mathbb{E}_Z (\ell(f, Z)) \\ &= R(f).\end{aligned}$$

The empirical risk is an **unbiased estimator** of the expected risk.

Finally, given \mathcal{D} , \mathcal{F} , and ℓ , “learning” aims at computing

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} \hat{R}(f; \mathcal{D}).$$

Finally, given \mathcal{D} , \mathcal{F} , and ℓ , “learning” aims at computing

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} \hat{R}(f; \mathcal{D}).$$

- Can we bound $R(f)$ with $\hat{R}(f; \mathcal{D})$?

Yes if f is not chosen using \mathcal{D} . Since the Z_n are independent, we just need to take into account the variance of $\hat{R}(f; \mathcal{D})$.

Finally, given \mathcal{D} , \mathcal{F} , and ℓ , “learning” aims at computing

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} \hat{R}(f; \mathcal{D}).$$

- Can we bound $R(f)$ with $\hat{R}(f; \mathcal{D})$?

Yes if f is not chosen using \mathcal{D} . Since the Z_n are independent, we just need to take into account the variance of $\hat{R}(f; \mathcal{D})$.

- Can we bound $R(f^*)$ with $\hat{R}(f^*; \mathcal{D})$?



Unfortunately not simply, and not without additional constraints on \mathcal{F} .

Finally, given \mathcal{D} , \mathcal{F} , and ℓ , “learning” aims at computing

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} \hat{R}(f; \mathcal{D}).$$

- Can we bound $R(f)$ with $\hat{R}(f; \mathcal{D})$?

Yes if f is not chosen using \mathcal{D} . Since the Z_n are independent, we just need to take into account the variance of $\hat{R}(f; \mathcal{D})$.

- Can we bound $R(f^*)$ with $\hat{R}(f^*; \mathcal{D})$?



Unfortunately not simply, and not without additional constraints on \mathcal{F} .

For instance if $|\mathcal{F}| = 1$, we can!

Note that in practice, we call “loss” both the functional

$$\ell : \mathcal{F} \times \mathcal{Z} \rightarrow \mathbb{R}$$

and the empirical risk minimized during training

$$\mathcal{L}(f) = \frac{1}{N} \sum_{n=1}^N \ell(f, z_n).$$

The end