



# Deep Learning for ECE

## EECE-580G

Introduction to Learning theory

# Notations

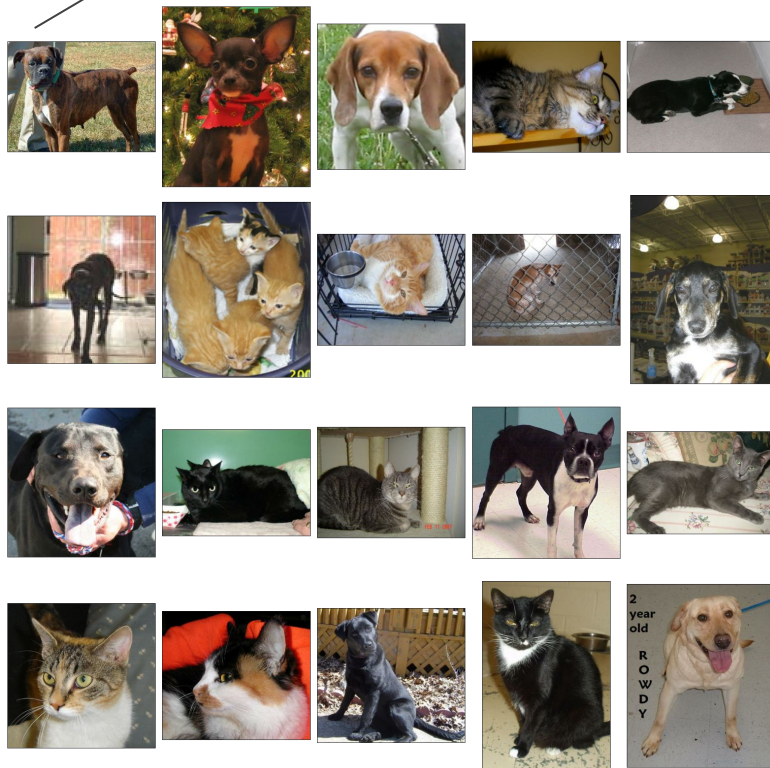
# Notations

— — —

- $m$  : samples or observations  $x^{(i)}$  are collected from a subject of interest (e.g. for a pet classification task, samples are photos from multiple pets)
- We consider these observations to be labeled, i.e. each sample  $x^{(i)}$  is assigned a label  $y^{(i)}$  (e.g. for each pet photo we collected, we know the breed of the pet appearing in it)
- This setting is called “supervised learning”<sup>1</sup>
- $D = (x^{(i)}, y^{(i)})_{i=1}^m$  : is the set of the available tuples of observations and labels, this is the dataset at hand
- $X \rightarrow x^{(i)}$ , and  $Y \rightarrow y^{(i)}$  :  $x^{(i)}$  and  $y^{(i)}$  are i.i.d. samples of the random variables  $X$  and  $Y$
- $X \sim P_X$ , and  $Y \sim P_Y$ ,
- $(X, Y) \rightarrow D$  the dataset is an observation of the joint random variable  $(X, Y)$
- $x^{(i)} \in \mathcal{X} = \mathbb{R}^n$  :  $n$  is the dimensionality of the data (e.g. considering the pet photos being of size  $H \times W$ ,  $n = 3 \cdot H \cdot W$ , the total number of pixels in one image represented in RGB space)
- $y^{(i)} \in \mathcal{Y} = \{0, 1\}$  : we consider the simple case of a binary label (e.g. only two pet breeds) the task we are performing is a binary classification task<sup>2</sup>
- $D \in \mathcal{D} = \mathbb{R}^n \times \{0, 1\}$ : the cartesian product of  $\mathcal{X}$  and  $\mathcal{Y}$

# Example

---  $x^{(0)}$

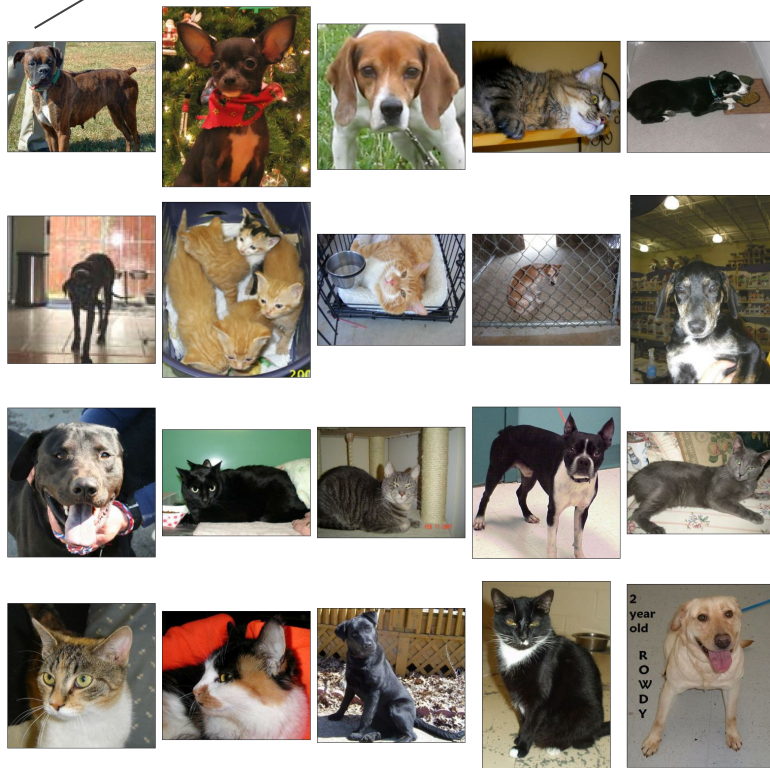


$y^{(0)}$

dog	dog	dog	cat	dog
dog	cat	cat	dog	dog
dog	cat	cat	dog	cat
cat	cat	dog	cat	dog

# Example

---  $x^{(0)}$

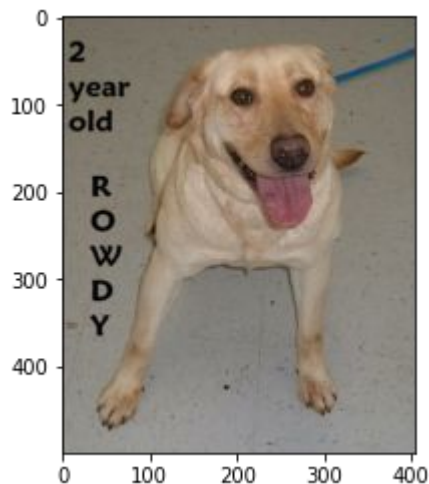


$y^{(0)}$

1	1	1	0	1
1	0	0	1	1
1	0	0	1	0
0	0	1	0	1

# Example

— — —



=

```
array([[108, 110, 113, ..., 110, 109],
       [110, 111, 113, ..., 108, 108],
       [112, 113, 114, ..., 108, 108],
       ...,
array([[102, 104, 107, ..., 106, 105, 105],
       [104, 105, 107, ..., 104, 104, 104],
       [106, 107, 108, ..., 104, 104, 104],
       ...,
array([[ 88,  90,  93, ...,  94,  93,  93],
       [ 90,  91,  93, ...,  92,  92,  92],
       [ 92,  93,  94, ...,  92,  92,  92],
       ...,
       [133, 137, 141, ..., 147, 144, 140],
       [135, 138, 142, ..., 148, 143, 138],
       [136, 139, 143, ..., 150, 148, 147]]], dtype=uint8)
```

# Notations

— — —

- We also consider that there exists a function  $f$ , that maps  $X$  to  $Y$ , i.e.  $f(X) = Y$ . This function is obviously unknown
- To link with the definition from the textbook:
  - Task  $T$  = binary classification task
  - Experiences  $E$  = the dataset  $D$
  - Performance measures  $P$  = loss function and notion of risk which will be introduced later

# Learning



# A learning algorithm

— — —

A learning algorithm:  $\mathcal{D} \rightarrow \mathcal{H} \subset \mathcal{Y}^{\mathcal{X}}$

$$D \mapsto \hat{f},$$

$\mathcal{H}$  = The hypothesis space, we can't look at the space of all functions ...

$\hat{f}$  = A function that “approximates” the ground truth function  
i.e.

$$\forall (x, y) \leftarrow (X, Y), \hat{f}(x) \approx f(x)$$

# 1-Nearest-Neighbor

— — —

$$\hat{f}(x) \mapsto y^{(j)}, j = \operatorname{argmin}_i \|x - x^{(i)}\|$$

- $\|\cdot\|$  can be any distance but we will use the Euclidean distance
- Maps  $x$  to the label of its closest neighbor in the dataset
- $\mathcal{H}$  ?

# 1-Nearest-Neighbor

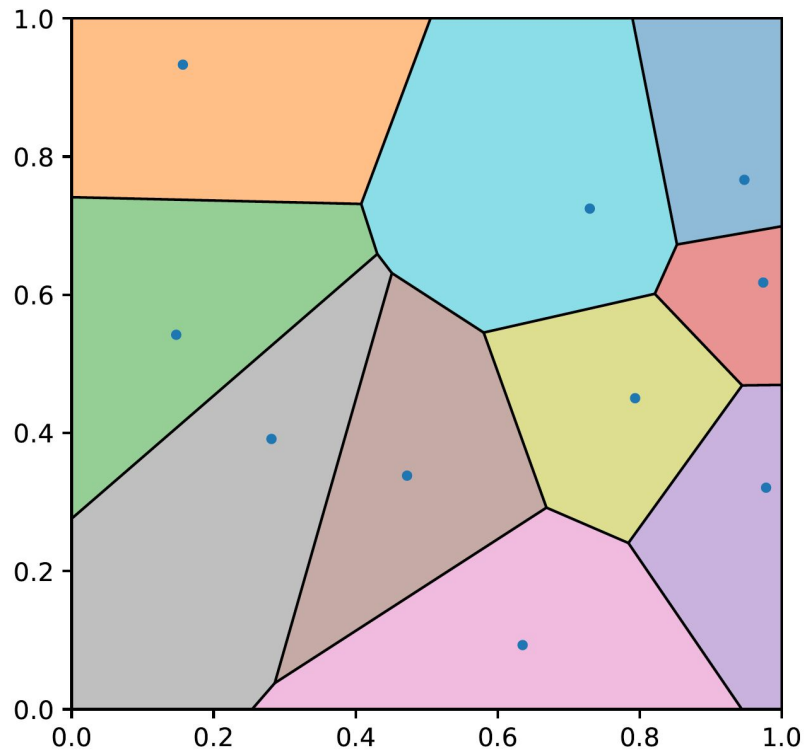
— — —

$$\hat{f}(x) \mapsto y^{(j)}, j = \operatorname{argmin}_i \|x - x^{(i)}\|$$

- $\| \cdot \|$  can be any distance but we will use the Euclidean distance
- Maps  $x$  to the label of its closest neighbor in the dataset
- $\mathcal{H}$  can be described using  $m$  parameters, the number of data points
- = the representational capacity  $C$

# 1-Nearest-Neighbor

— — —



# Empirical Risk Minimization

# Risk - loss function $l(\hat{f}(x), y)$

— — —

- We need to quantify “goodness” of fit of the approximation
  - 0-1 loss
  - Weighted 0-1 loss with a cost matrix
  - “Smooth” surrogate losses (later in the course)
  - ...

# Empirical risk - Generalization risk

— — —

$$\tilde{L}(h) = \frac{1}{m} \sum_{i=1}^m l(h(x^{(i)}), y^{(i)})$$

- We don't know the distributions of  $(X, Y)$
- We are interested in the “generalization risk”

$$L(h) = E [l(h(X), Y)]$$

- Luckily when samples are i.i.d.

$$E [\tilde{L}(h)] = E \left[ \frac{1}{m} \sum_{i=1}^m l(h(x^{(i)}), y^{(i)}) \right] = E [l(h(X), Y)] = L(h)$$

# Empirical risk minimizer

— — —

$$\tilde{L}(h) = \frac{1}{m} \sum_{i=1}^m l(h(x^{(i)}), y^{(i)})$$

- So let's pick the approximation function as:

$$\hat{f} = \operatorname{argmin}_{h \in \mathcal{H}} \tilde{L}(h)$$

- The “empirical risk minimizer”
- How does the empirical risk of empirical risk minimizer link to its generalization risk

$$E \left[ \tilde{L}(\hat{f}) \right] \stackrel{?}{=} L(\hat{f})$$



# Empirical risk minimizer

— — —

$$\tilde{L}(h) = \frac{1}{m} \sum_{i=1}^m l(h(x^{(i)}), y^{(i)})$$

- So let's pick the approximation function as:

$$\hat{f} = \operatorname{argmin}_{h \in \mathcal{H}} \tilde{L}(h)$$

- The “empirical risk minimizer”
- How does the empirical risk of empirical risk minimizer link to its generalization risk

$$E \left[ \tilde{L}(\hat{f}) \right] \neq L(\hat{f})$$

- Why?

# Empirical risk minimizer

— — —

- $l(\hat{f}(x^{(i)}), y^{(i)})$  Are not necessarily i.i.d anymore
- $\hat{f}$  Is dependent on the entire dataset!  $\hat{f} = \operatorname{argmin}_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m l(h(x^{(i)}), y^{(i)})$
- In general, the empirical risk of empirical risk minimizer is a **biased estimator** of its generalization risk

$$E \left[ \tilde{L}(\hat{f}) \right] \neq L(\hat{f})$$

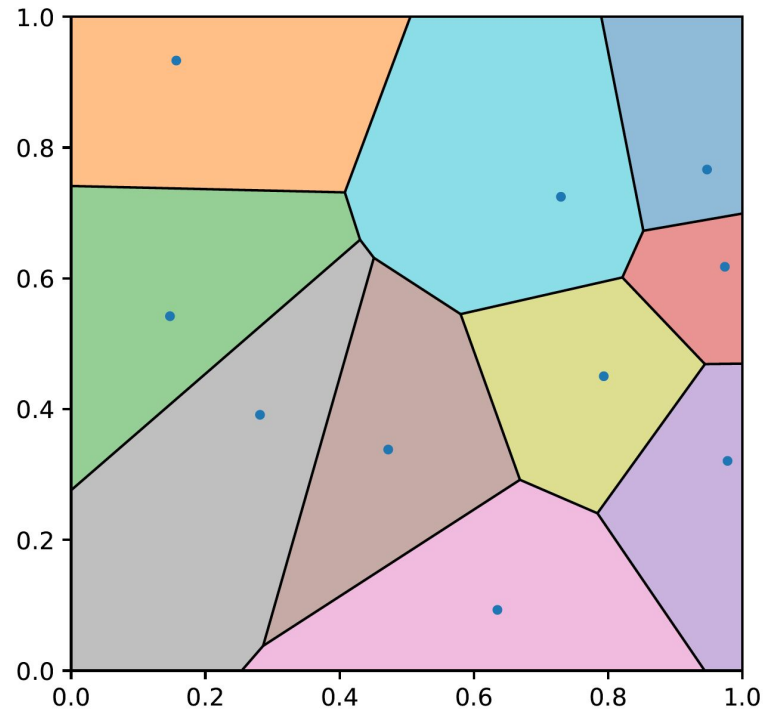
# Example - 1-Nearest-Neighbor

---

- $\tilde{L}(\hat{f}) = 0$

Because the 1-Nearest Neighbor to each datapoint is itself!

- But in almost all cases,  $L(\hat{f}) \neq 0$



# Risk decomposition

# Risk decomposition

---

$$\hat{f} = \operatorname{argmin}_{h \in \mathcal{H}} \tilde{L}(h) \quad \text{Empirical risk minimizer}$$

$$f^* = \operatorname{argmin}_{h \in \mathcal{H}} L(h) \quad \text{Generalization risk minimizer}$$

- $f^*$  Is unknown and impossible to compute, we can't compute the expectations...
- Is  $L(\hat{f})$  close enough to  $L(f^*)$  ?

# Risk decomposition

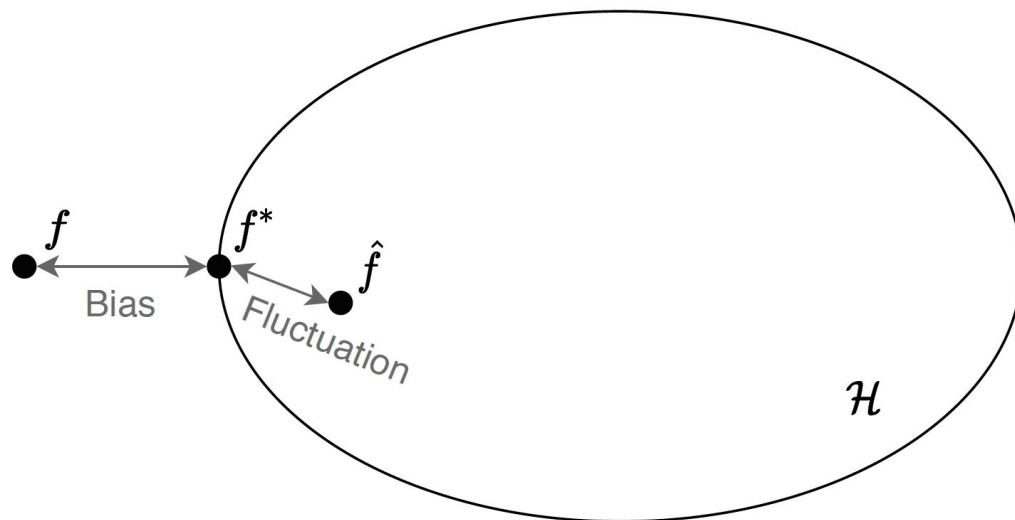
— — —

$$L(f^*) \leq L(\hat{f}) \leq \underbrace{L(f^*)}_{\text{Bias}} + \underbrace{2 \max_{h \in \mathcal{H}} |L(h) - \tilde{L}(h)|}_{\text{Fluctuation (variance)}}$$

- Proof in handout

# Risk decomposition

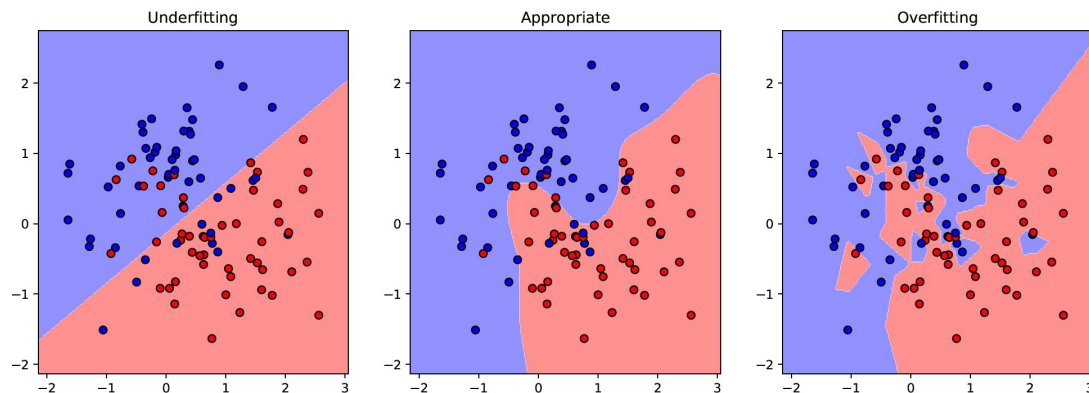
— — —



# Bias variance trade-off

— — —

- Larger  $\mathcal{H} \longrightarrow$  **smaller bias**, but **higher fluctuations**
- The choice of  $\mathcal{H}$  is not trivial, needs to be validated by the data too
- This is called the **bias-variance trade-off**





# Training/validation/testing sets

# The need for a Validation set

- — —
- Remember that  $E \left[ \tilde{L}(\hat{f}) \right] \neq L(\hat{f})$
  - But what if we computed  $\tilde{L}(\hat{f})$  using a different set of data, not used to find  $\hat{f}$
  - I.e. we split  $D$  into  $D^{\text{Train}}$  and  $D^{\text{Validation}}$

$$\tilde{L}^{\text{Validation}}(\hat{f}) = \frac{1}{m} \sum_{(x^{(i)}, y^{(i)}) \in D^{\text{Validation}}} l(\hat{f}(x^{(i)}), y^{(i)})$$

- Is unbiased now!

# The need for a Validation set

— — —

The diagram illustrates the partitioning of a dataset. A large light green rectangle at the top represents the 'Entire Dataset'. Below it, the dataset is split into two horizontal sections. The left section, also light green, represents the 'Training set, used for determining  $\hat{f}$ '. The right section, light blue, represents the 'Validation set, used for approximating  $L(\hat{f})$ '. The training set section is wider than the validation set section.

Entire Dataset

Training set, used for determining  $\hat{f}$

Validation set, used for approximating  $L(\hat{f})$

# The need for a Validation set

— — —

Entire Dataset

Training set, used for determining  $\hat{f}$

Validation set, used for approximating  $L(\hat{f})$

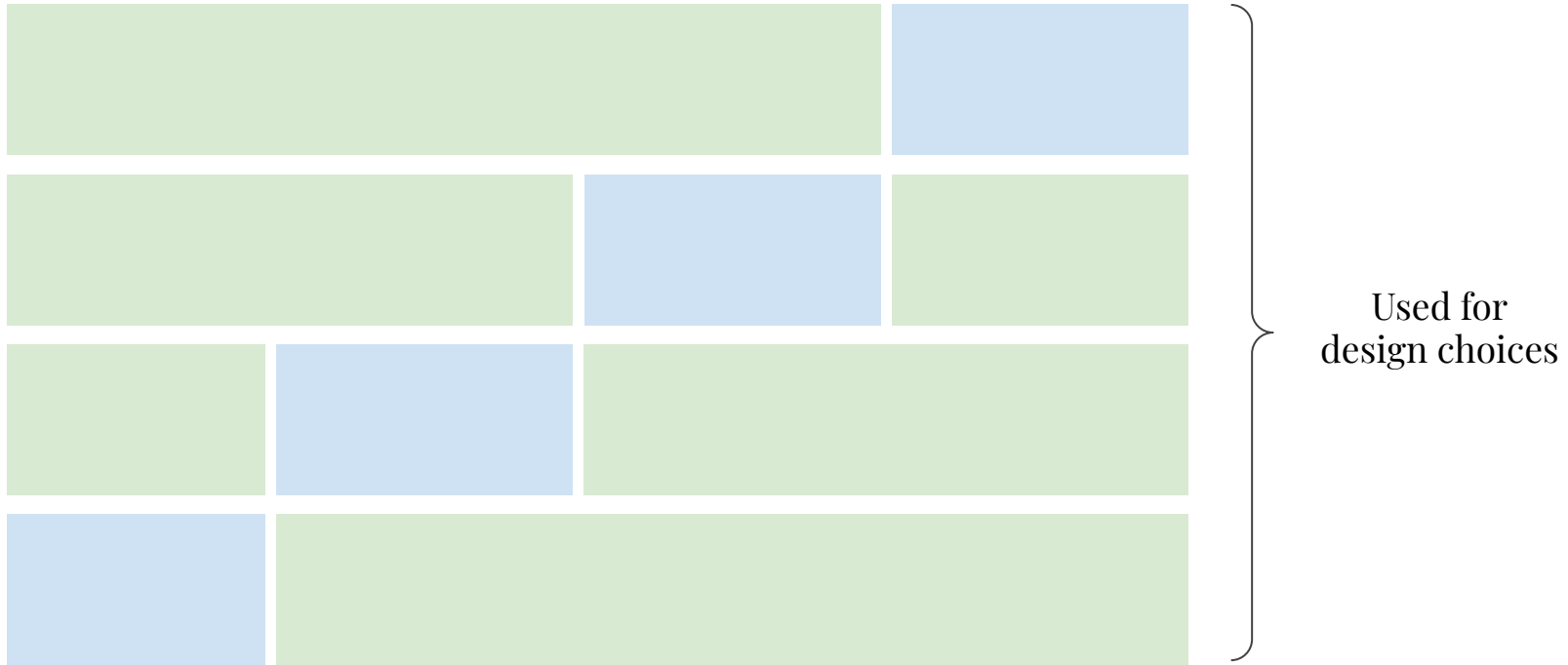
Train

Validation used for making design choices

Test set used for approximating  $L(\hat{f})$

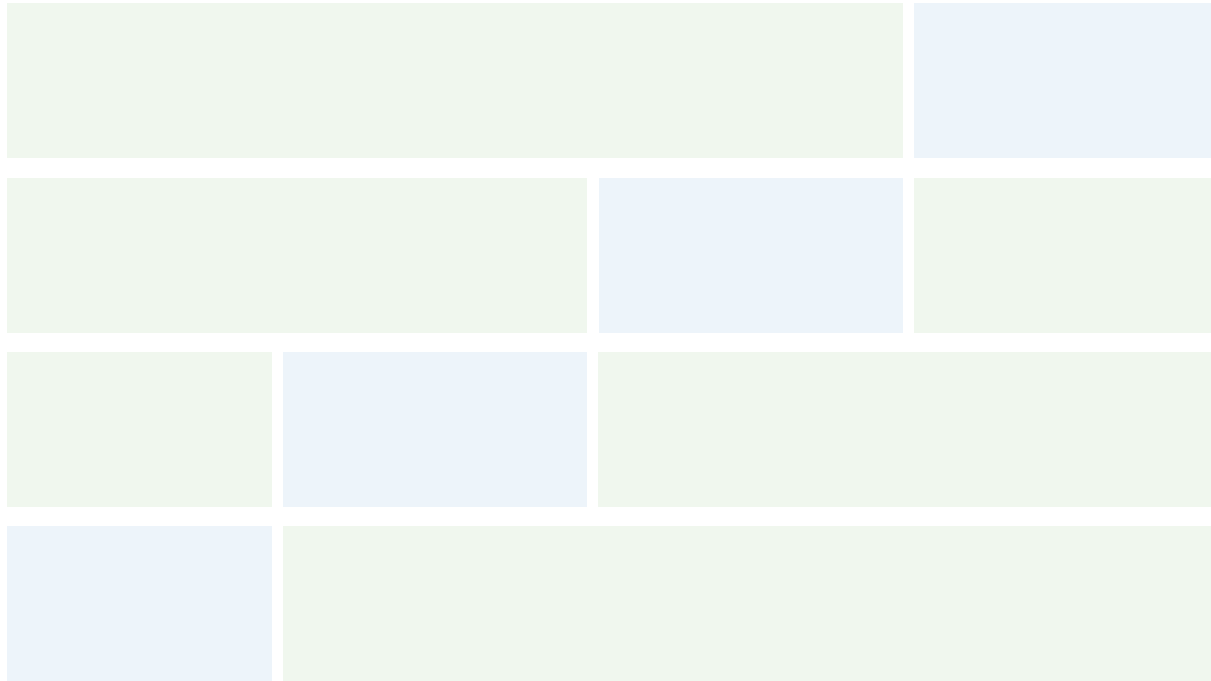
# The need for a Validation set

— — —



# The need for a Validation set

— — —



Test set used for approximating  $L(\hat{f})$

**End**