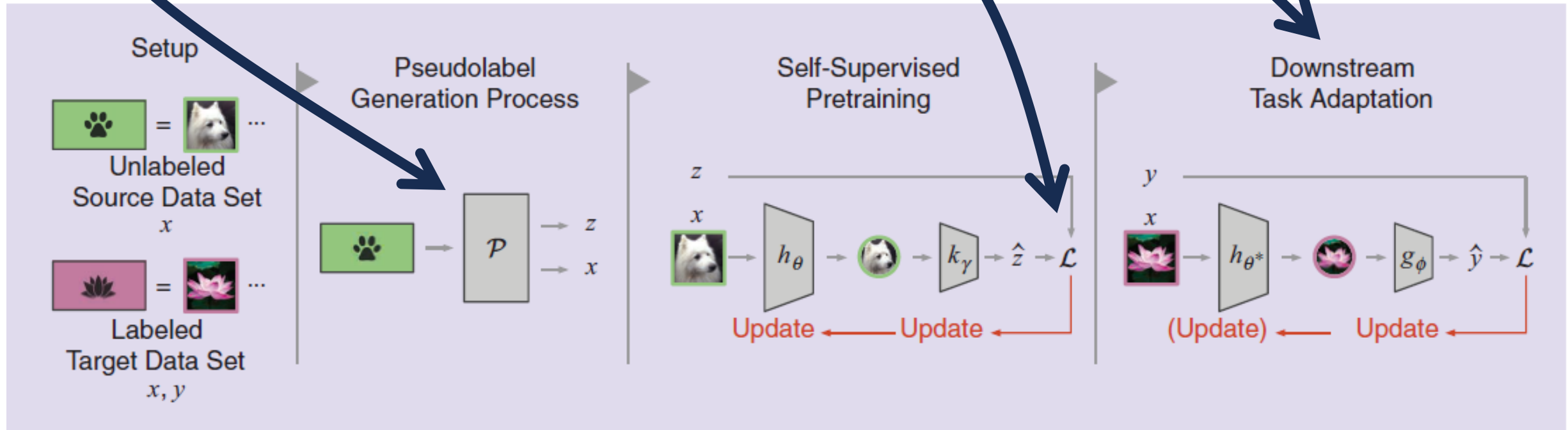**FIGURE 1.** Contrasting supervised, unsupervised and self-supervised learning paradigms for training a model $f$ using raw data $x$, labels $y$, and loss function $\mathcal{L}$. Self-supervision methods introduce pretext tasks $\mathcal{P}$ that generate pseudolabels $z$ for discriminative training of $f$.

# Module – 2

Basics of Foundation Models, SSL formalisms, definitions and examples of

- **Pretext tasks**
- **Losses**
- **Downstream adaptations**

in a domain-agnostic setting



**FIGURE 2.** The self-supervised workflow starts with an unlabeled source data set and a labeled target data set. As defined by the pretext task, pseudolabels are programmatically generated from the unlabeled set. The resulting inputs, $x$ and pseudolabels $z$, are used to pretrain the model $k_\gamma(h_\theta(\cdot))$—composed of feature extractor $h_\theta$ and output $k_\gamma$ modules—to solve the pretext task. After pretraining is complete, the learned weights $\theta^*$ of the feature extractor $h_{\theta^*}$ are transferred and used together with a new output module $g_\phi$ to solve the downstream target task.

## 4.10    Theory

*Authors: Aditi Raghunathan, Sang Michael Xie, Ananya Kumar, Niladri Chatterji, Rohan Taori, Tatsunori Hashimoto, Tengyu Ma*
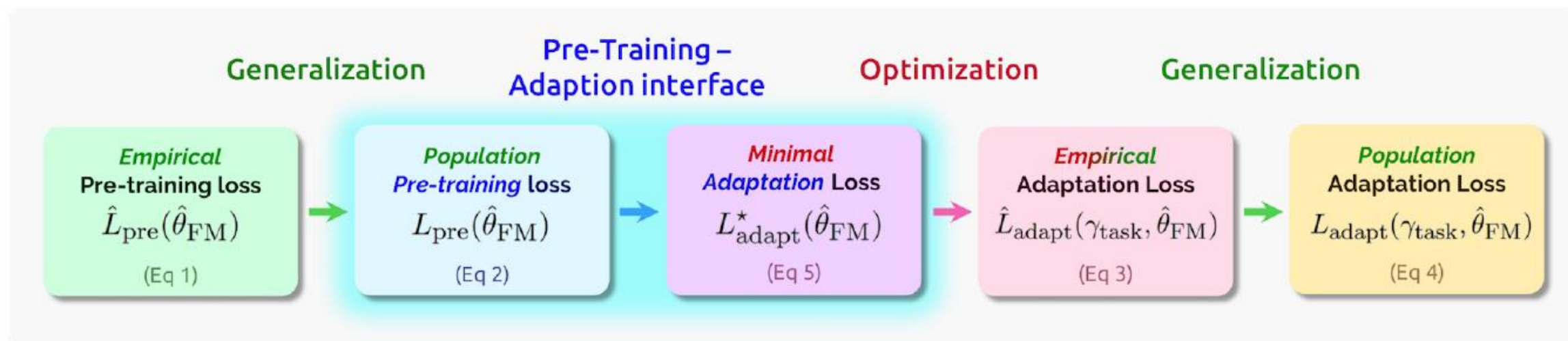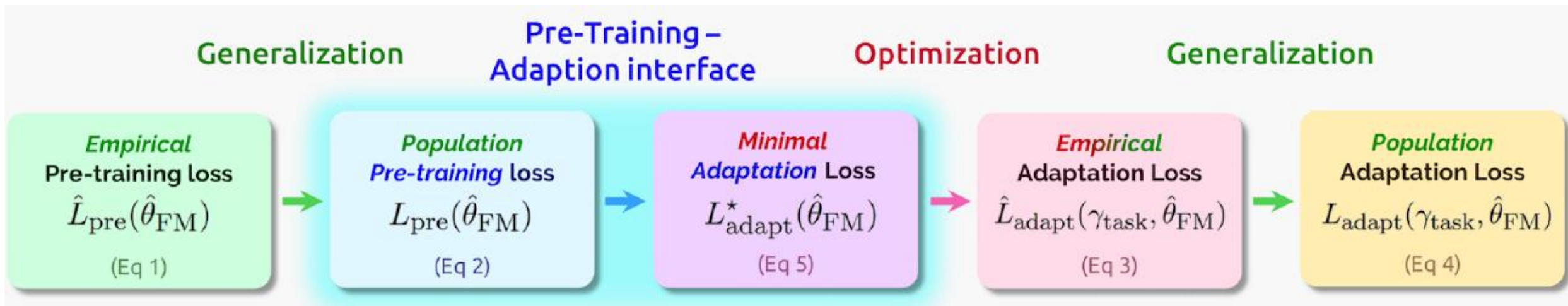


Fig. 22.   The analysis of foundation models from pretraining on diverse data to downstream performance on adapted tasks involves capturing the relation between different loss terms as shown above. The main challenge is to analyze the highlighted pretraining-adaptation interface which requires reasoning carefully about the population losses in addition to the model architecture, losses and data distributions of the pretraining and adaptation stages (§4.10.2: THEORY-INTERFACE). Analysis of generalization and optimization largely reduces to their analysis in standard supervised learning.
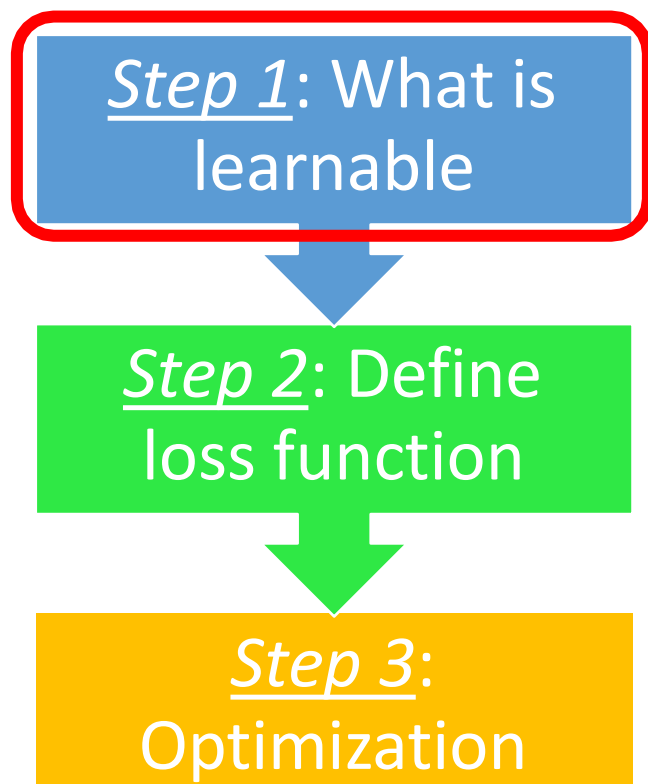
3

| Generalization | Pre-Training – Adaption interface | Optimization | Generalization |
|---|---|---|---|

**Empirical Pre-training loss**
$\hat{L}_{\mathrm{pre}}(\hat{\theta}_{\mathrm{FM}})$
(Eq 1)

**Population Pre-training loss**
$L_{\mathrm{pre}}(\hat{\theta}_{\mathrm{FM}})$
(Eq 2)

**Minimal Adaptation Loss**
$L^{\star}_{\mathrm{adapt}}(\hat{\theta}_{\mathrm{FM}})$
(Eq 5)

**Empirical Adaptation Loss**
$\hat{L}_{\mathrm{adapt}}(\gamma_{\mathrm{task}}, \hat{\theta}_{\mathrm{FM}})$
(Eq 3)

**Population Adaptation Loss**
$L_{\mathrm{adapt}}(\gamma_{\mathrm{task}}, \hat{\theta}_{\mathrm{FM}})$
(Eq 4)

As shown in Figure 22, the main missing link beyond standard supervised theory is:

*Under what conditions does a small population pretraining loss $L_{pre}(\hat{\theta}_{\mathrm{FM}})$ imply a small minimal adaptation loss $L^{\star}_{adapt}(\hat{\theta}_{\mathrm{FM}})$ and why?*
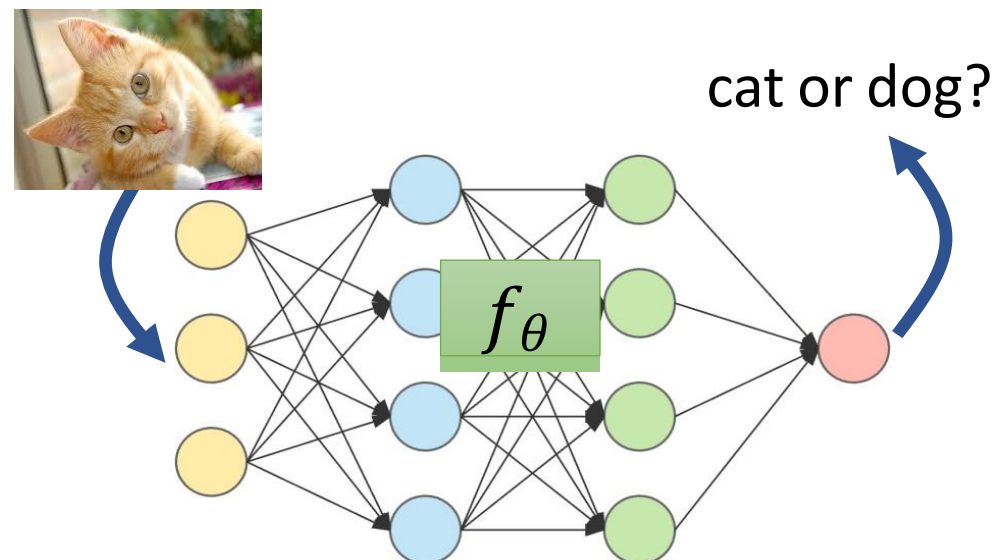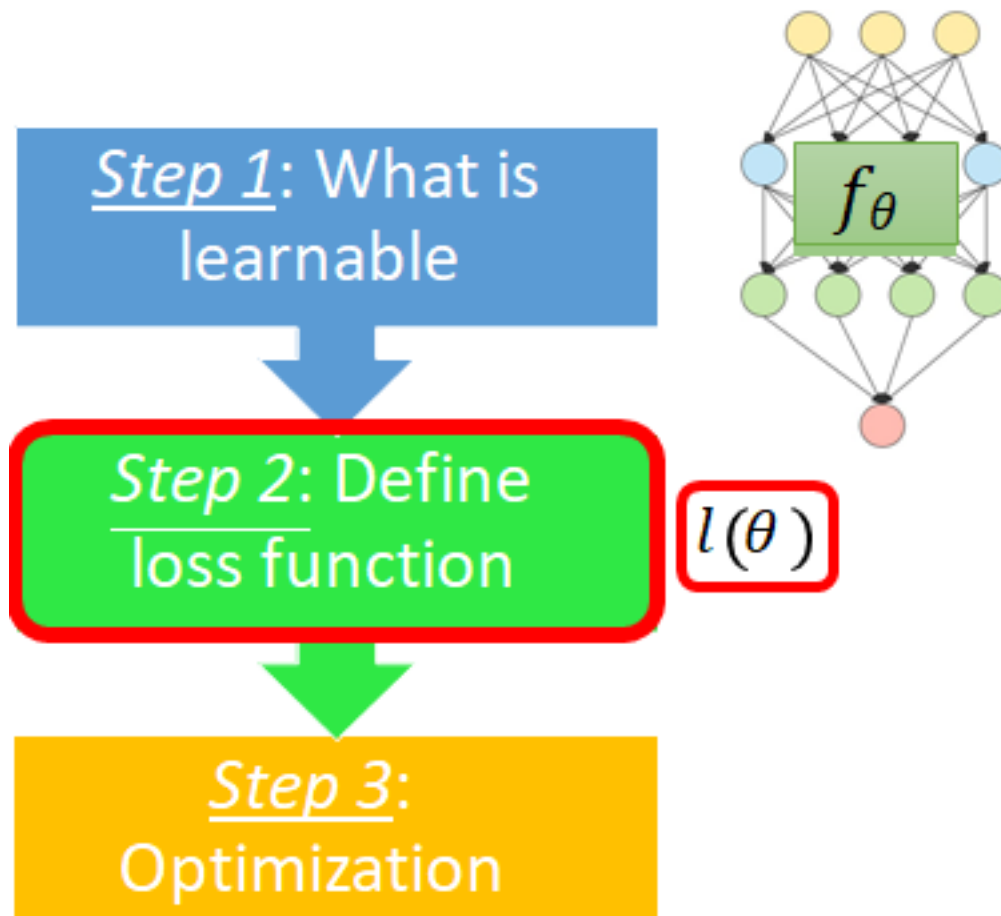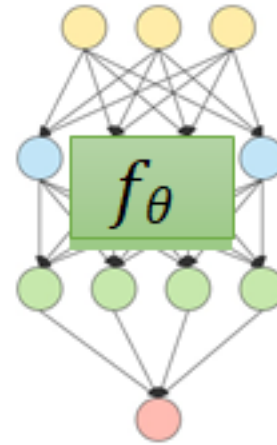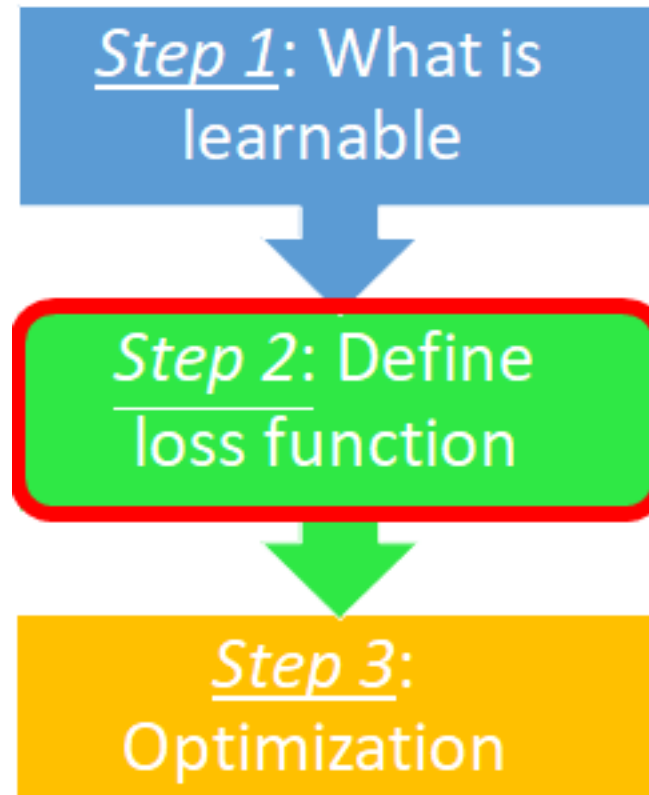
# Deep Learning pipeline

Dog-Cat Classification

$f\,($  $) = $ "cat"

**Step 1**: What is learnable

**Step 2**: Define loss function

**Step 3**: Optimization



cat or dog?

$f_\theta$

Weights and biases of neurons are learnable.

Using $\theta$ to represent the learnable parameters.

**Step 1:** What is learnable

$f_\theta$

**Step 2:** Define loss function

$l(\theta)$

**Step 3:** Optimization

6

Training Examples

Step 1: What is learnable

$f_\theta$

Step 2: Define loss function

$l(\theta)$

Step 3: Optimization

cat    dog         cat    dog

Ground Truth

7

Step 1: What is learnable

Step 2: Define loss function

$l(\theta)$

Step 3: Optimization

Training Examples

$f_\theta$

$f_\theta$

$f_\theta$

cat    dog

cat    dog

cat    dog

cat    dog

Ground Truth

*Training Examples*

Step 1: What is learnable

Step 2: Define loss function $l(\theta)$

Step 3: Optimization

cat   dog

cat   dog

Cross-entropy $d_1$   $d_2$

cat   dog

cat   dog

*Ground Truth*

$$l(\theta) = \sum_{k=1}^{K} d_k$$

Step 1: What is learnable

Step 2: Define loss function

Step 3: Optimization
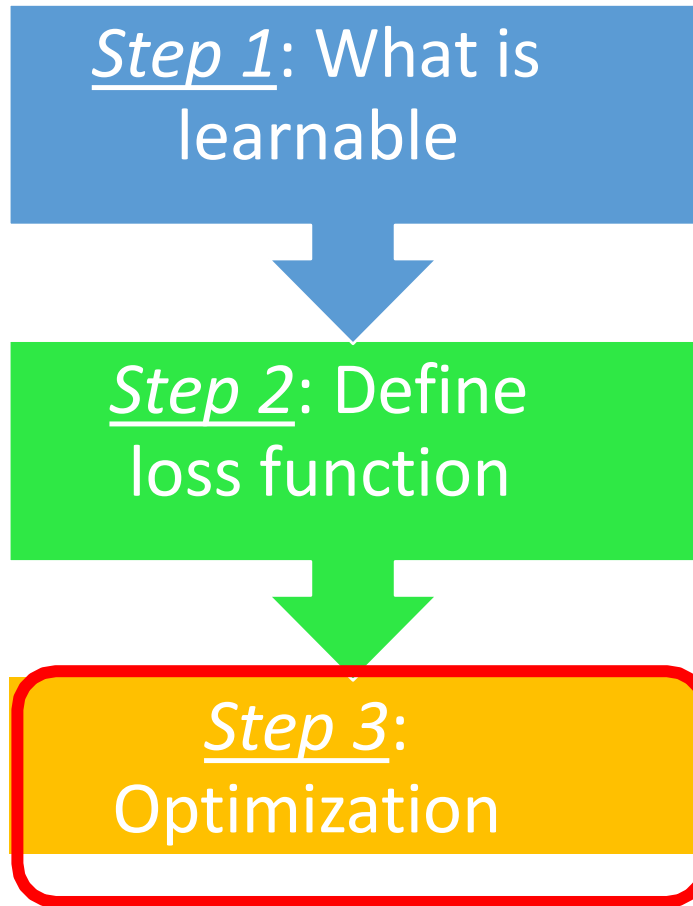
loss: $l(\theta) = \sum_{k=1}^{K} d_k$   sum over training examples

$\hat{\theta} = arg \min_{\theta} l(\theta)$
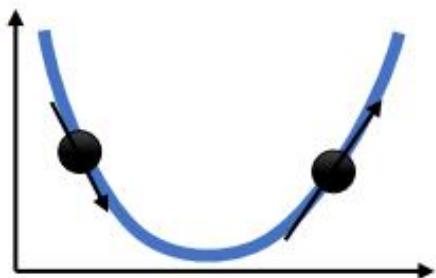
done by gradient descent

$f_{\hat{\theta}}$ is the function learned by learning algorithm from data

# Gradient descent

An algorithm:

1. Find a *direction* $v$ where $\mathcal{L}(\theta)$ decreases
2. $\theta \leftarrow \theta + \alpha v$

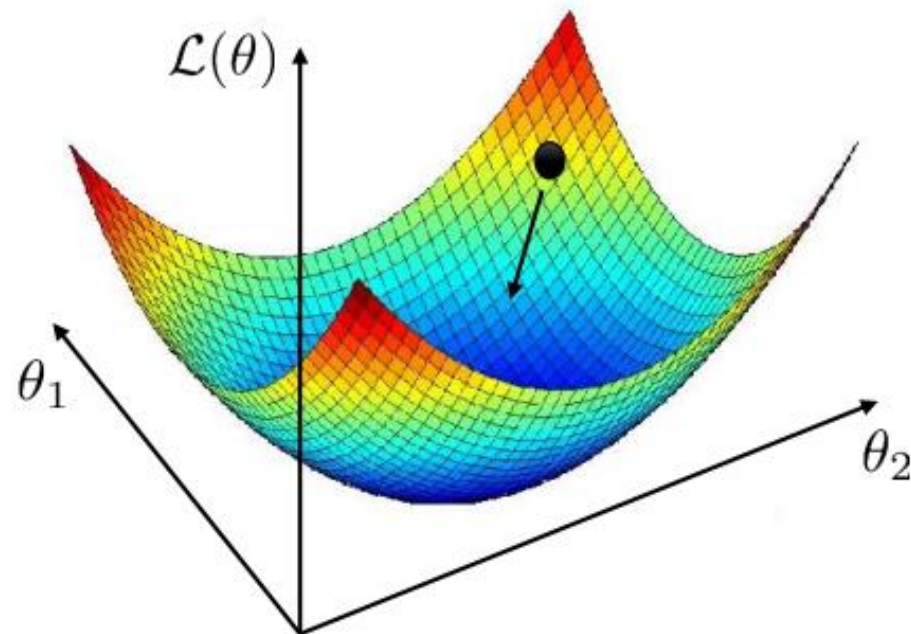Which way does $\mathcal{L}(\theta)$ decrease?

negative slope = go to the right

positive slope = go to the left

in general:

for each dimension, go in the direction opposite the slope **along that dimension**

$$v_1 = -\frac{d\mathcal{L}(\theta)}{d\theta_1} \quad v_2 = -\frac{d\mathcal{L}(\theta)}{d\theta_2} \quad \text{etc.}$$

gradient:

$$\nabla_\theta \mathcal{L}(\theta) = \begin{pmatrix} \dfrac{d\mathcal{L}(\theta)}{d\theta_1} \\[2mm] \dfrac{d\mathcal{L}(\theta)}{d\theta_2} \\[1mm] \vdots \\[1mm] \dfrac{d\mathcal{L}(\theta)}{d\theta_n} \end{pmatrix}$$
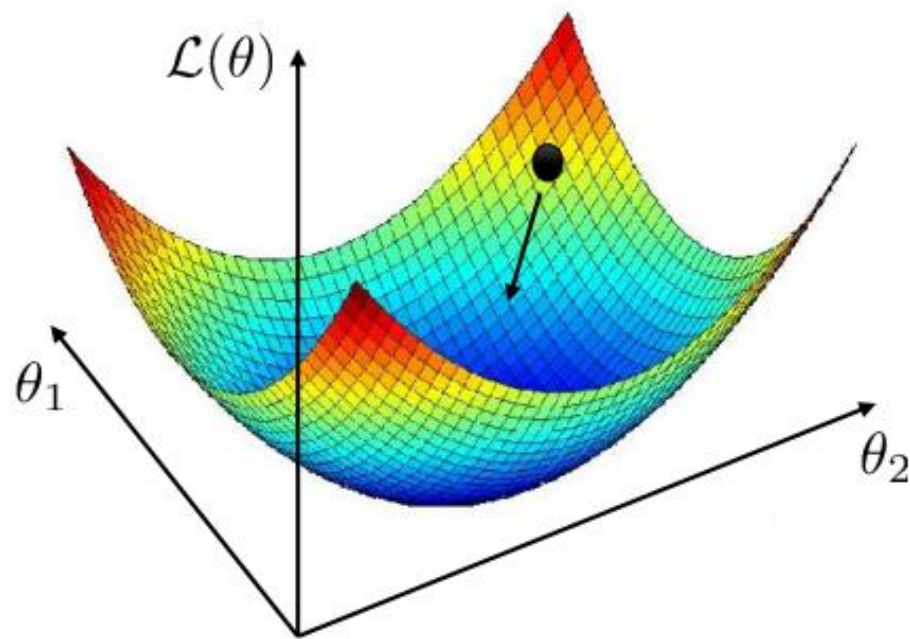
# Gradient descent

An algorithm:

1. Find a *direction* $v$ where $\mathcal{L}(\theta)$ decreases
2. $\theta \leftarrow \theta + \alpha v$

Gradient descent:

1. Compute $\nabla_\theta \mathcal{L}(\theta)$
2. $\theta \leftarrow \theta - \alpha \nabla_\theta \mathcal{L}(\theta)$



$$\nabla_\theta \mathcal{L}(\theta) = \begin{pmatrix} \dfrac{d\mathcal{L}(\theta)}{d\theta_1} \\ \dfrac{d\mathcal{L}(\theta)}{d\theta_2} \\ \vdots \\ \dfrac{d\mathcal{L}(\theta)}{d\theta_n} \end{pmatrix}$$

Thank you !!