CS-C3240 - Machine Learning

# Model Regularization

Data Augmentation. Soft Model-Selection. Transfer Learning. Multi-Task Learning. Semi-Supervised Learning.

Alexander Jung

# What I want to teach you today:

- basic idea of regularization

- regularization as soft model selection

- basic idea of data augmentation

- equivalence between regularization and data aug.

# What is ML ?

**informal:** learn hypothesis out of a hypothesis space or "model" that incurs minimum loss when predicting labels of datapoints based on their features

"training error"

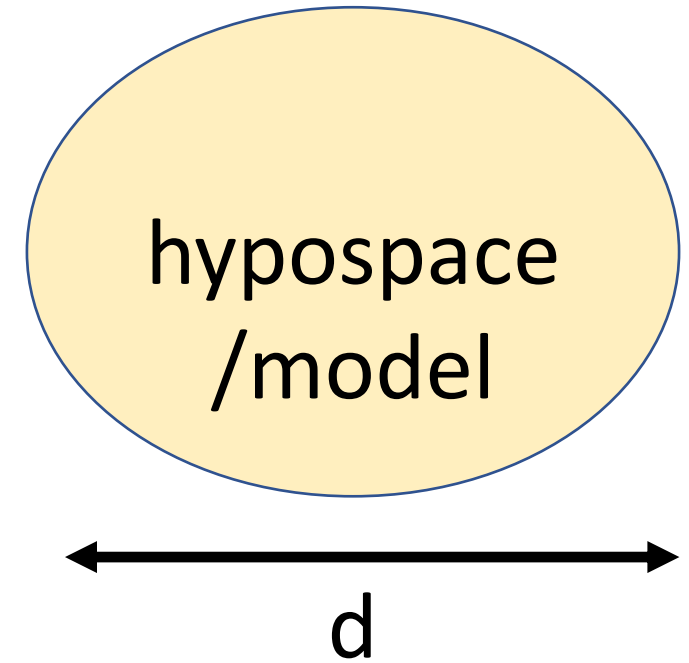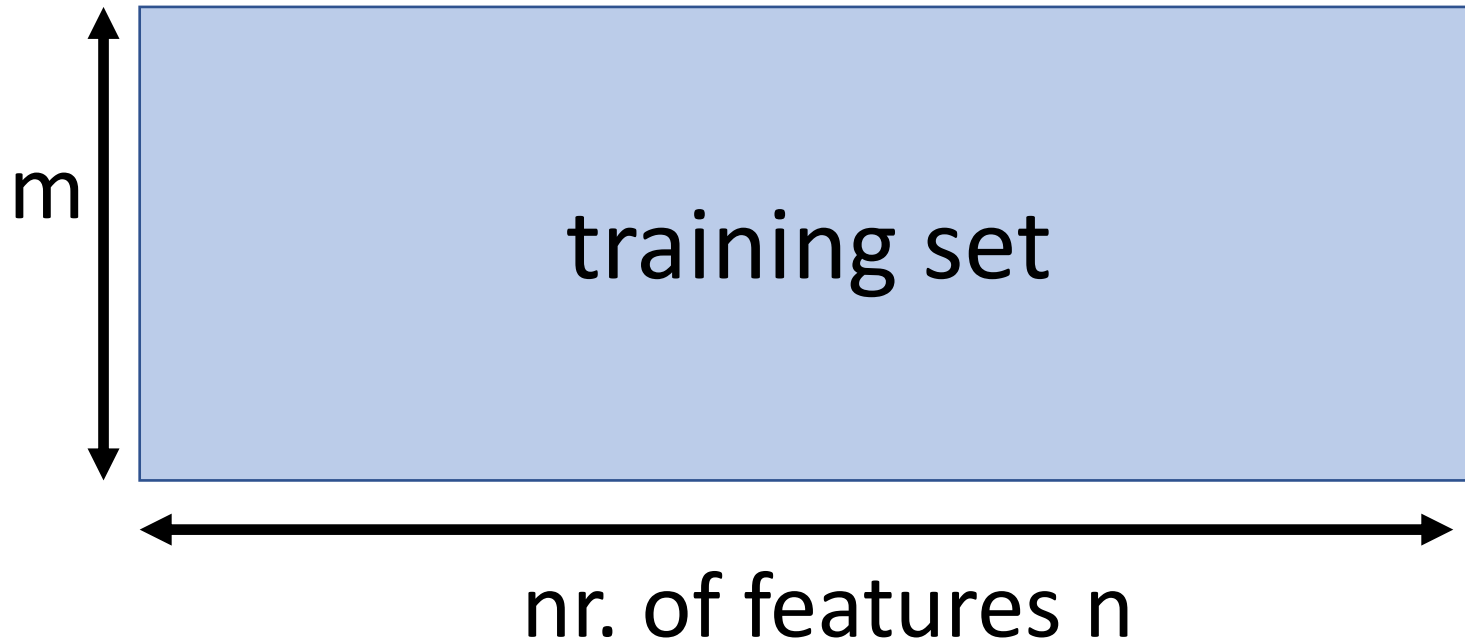$$\hat{h} = \operatorname*{argmin}_{h \in \mathcal{H}} \mathcal{E}(h|\mathcal{D})$$

$$\stackrel{(2.12)}{=} \operatorname*{argmin}_{h \in \mathcal{H}} (1/m) \sum_{i=1}^{m} \mathcal{L}((\mathbf{x}^{(i)}, y^{(i)}), h).$$
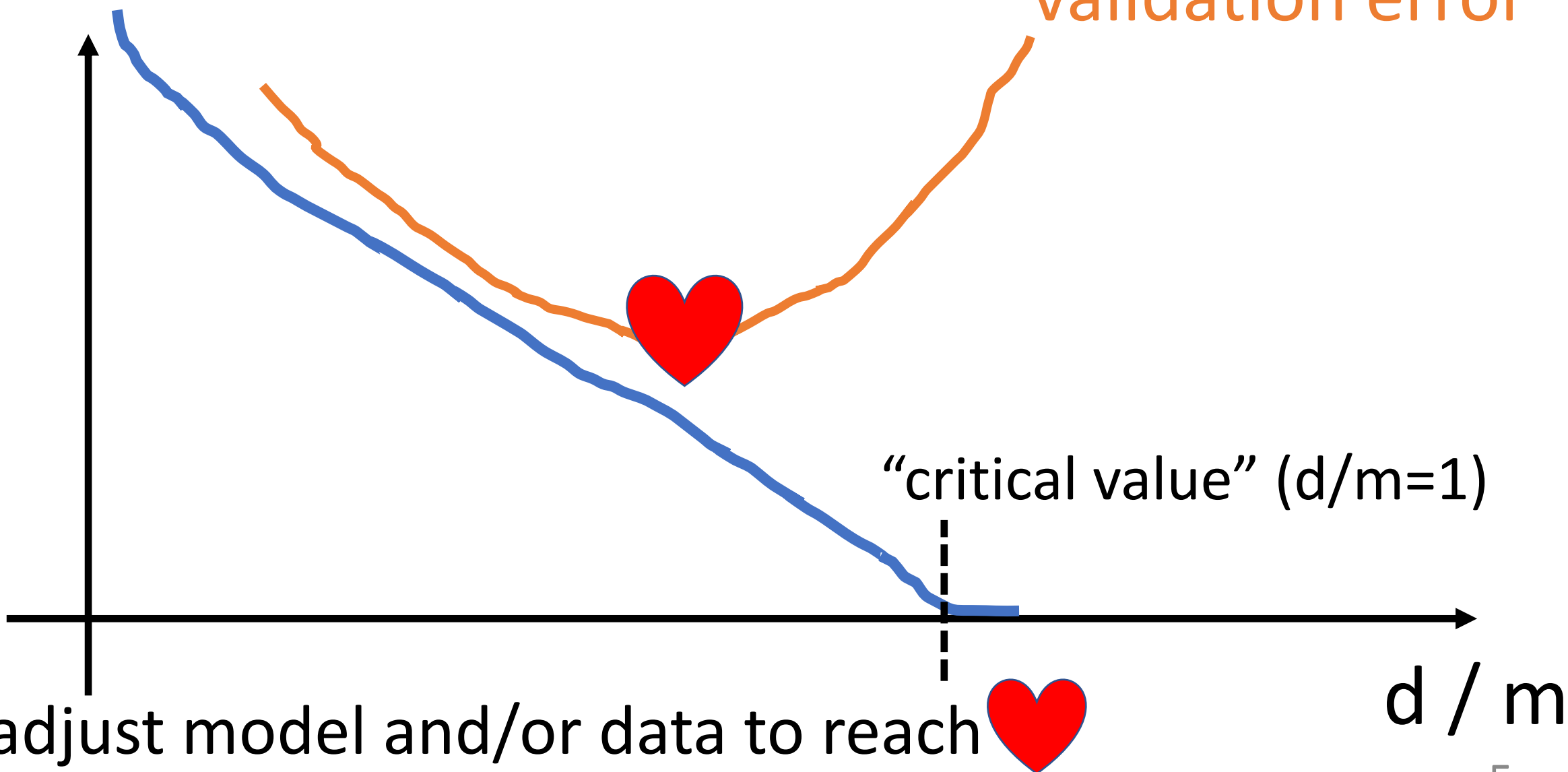
see Ch. 4.1 of mlbook.cs.aalto.fi

10.3.2021

3

# Data and Model Size



training set

m

nr. of features n

hypospace /model

d

crucial parameter is the ratio d/m

training error

validation error

"critical value" (d/m=1)

d / m

adjust model and/or data to reach ❤

bring d/m below critical value 1:
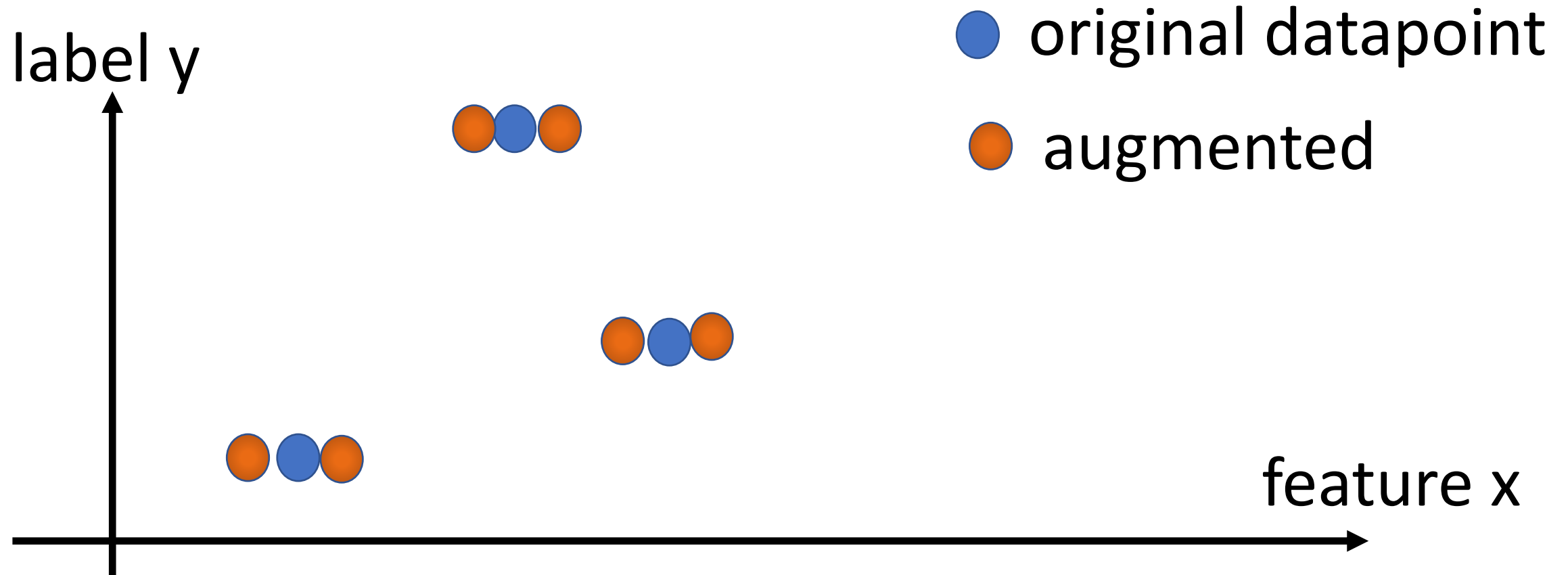
- increase m by using more training data

- decrease d by using smaller hypothesis space

bring d/m below critical value 1:

- <span style="color:red">increase m by using more training data</span>

- decrease d by using smaller hypothesis space

# Data Augmentation

# add a bit of noise to features



label y

original datapoint

augmented

feature x

we have increased the dataset by factor 3 !

# rotated cat image is still cat image

# flipped cat image is still cat image

# shifted cat image is still cat image

bring d/m below critical value 1:

- increase m by using more training data

- decrease d by using smaller hypothesis space

replace original ERM
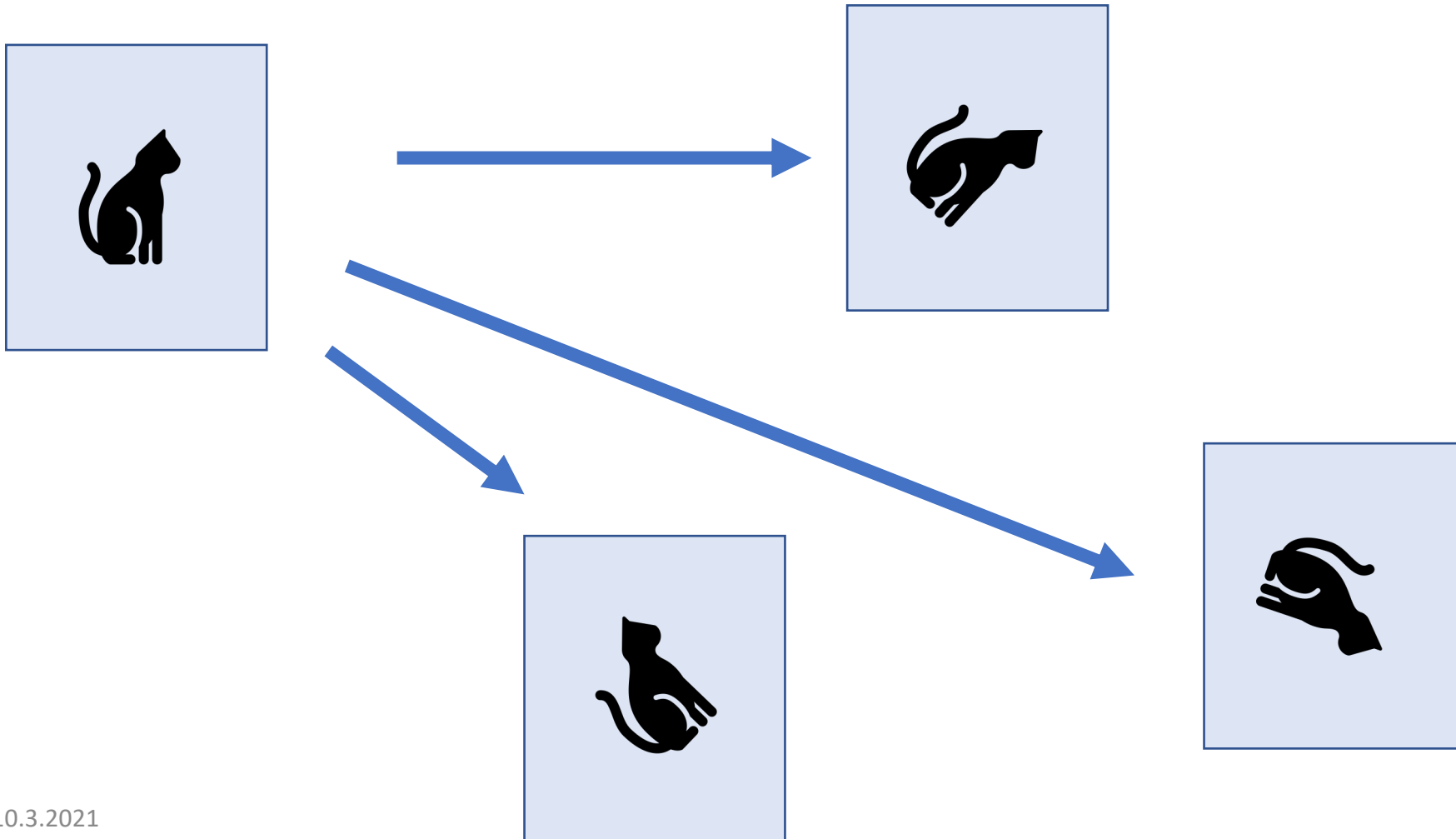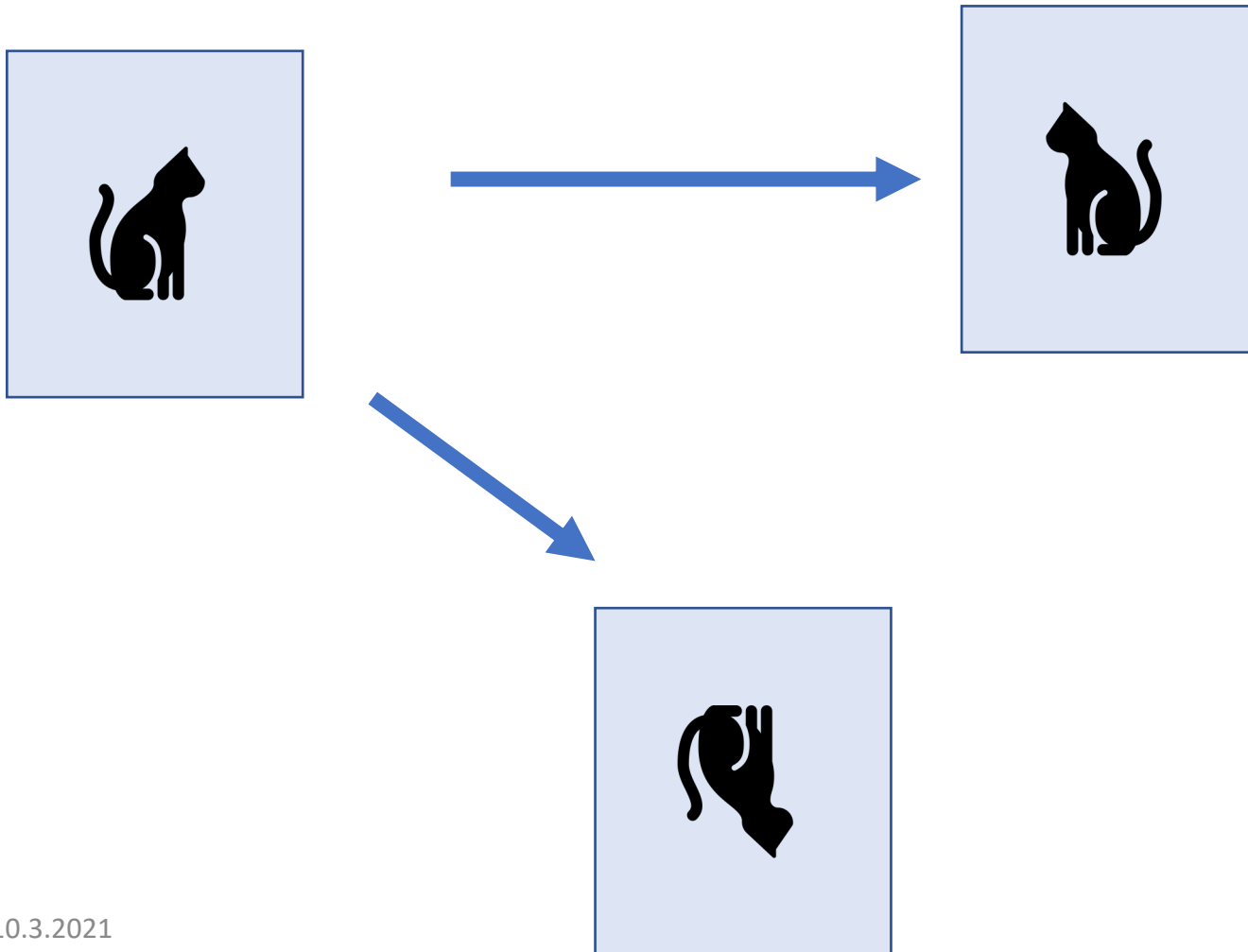
$$\min_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}\big((x^{(i)}, y^{(i)}), h\big)$$

with ERM on smaller $\widehat{\mathcal{H}} \subset \mathcal{H}$

$$\min_{h \in \widehat{\mathcal{H}}} \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}\big((x^{(i)}, y^{(i)}), h\big)$$

# Nested Models



degree 1 polyn.

degree 2 polyn.

degree 3 polyn.

# Prune Hypospace by Early Stopping



10 iterations

100 iterations

10000 iterations

# Soft Model Pruning via Regularization

# Regularized ERM

learn hypothesis $h$ out of
model (hypospace) $\mathcal{H}$ by minimizing

$$\underbrace{\frac{1}{m}\sum_{i=1}^{m}\mathcal{L}\big((x^{(i)}, y^{(i)}), h\big)}_{\text{average loss on training set (empirical risk of h)}} + \underbrace{\lambda\mathcal{R}(h)}_{\text{loss increase for datapoints outside training set}}$$

# Regularized Linear Regression

- squared error loss

- linear hypothesis map $h(x) = w^T x = w_1 x_1 + \cdots + w_n x_n$

$$\frac{1}{m} \sum_{i=1}^{m} \left( y^{(i)} - w^T x^{(i)} \right)^2 + \lambda \mathcal{R}(w)$$

- ridge regression uses $\mathcal{R}(w) = \|w\|_2^2 = w_1^2 + \cdots + w_n^2$

- Lasso uses $\mathcal{R}(w) = \|w\|_1 = |w_1| + \cdots + |w_n|$

# Regularization = Implicit Pruning!

$$\min_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}\big((x^{(i)}, y^{(i)}), h\big) + \lambda \mathcal{R}(h)$$

equivalent to

$$\min_{h \in \mathcal{H}^{(\lambda)}} \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}\big((x^{(i)}, y^{(i)}), h\big)$$

with pruned model $\mathcal{H}^{(\lambda)} \subset \mathcal{H}$

# Regularization = "Soft" Model Selection



$\lambda_1 \quad < \quad \lambda_2 \quad < \quad \lambda_3$

$\mathcal{H}^{(\lambda_1)}$

$\mathcal{H}^{(\lambda_2)}$

$\mathcal{H}^{(\lambda_3)}$

$\lambda$

# Regularization does implicit Data Augmentation

# augment with (infinitely many) realizations of RV!

label y

○ original datapoint

○ augmented

○ = ● + "noise"

feature x

# Regularization =Implicit Data Aug.



label y

h(x)

● raw datapoint

● "perturbed" datapoint

$$\frac{1}{m}\sum_{i=1}^{m}\mathcal{L}\big((x^{(i)}, y^{(i)}), h\big) + \lambda\mathcal{R}(h)$$

see Chapter 7.3 of mlbook.cs.aalto.fi

feature x

# Transfer Learning via Regularization

- Problem I: classify image as "shows border collie" vs. "not"

- Problem II: classify image as "shows a dog" vs. "not"

- ML Problem I is our main interest

- only little training data $\mathcal{D}^{(1)}$ for Problem I

- much more labeled data $\mathcal{D}^{(2)}$ for Problem II

- pre-train a hypothesis on $\mathcal{D}^{(2)}$ , fine-tune on $\mathcal{D}^{(1)}$

The following 81 files are in this category, out of 81 total.
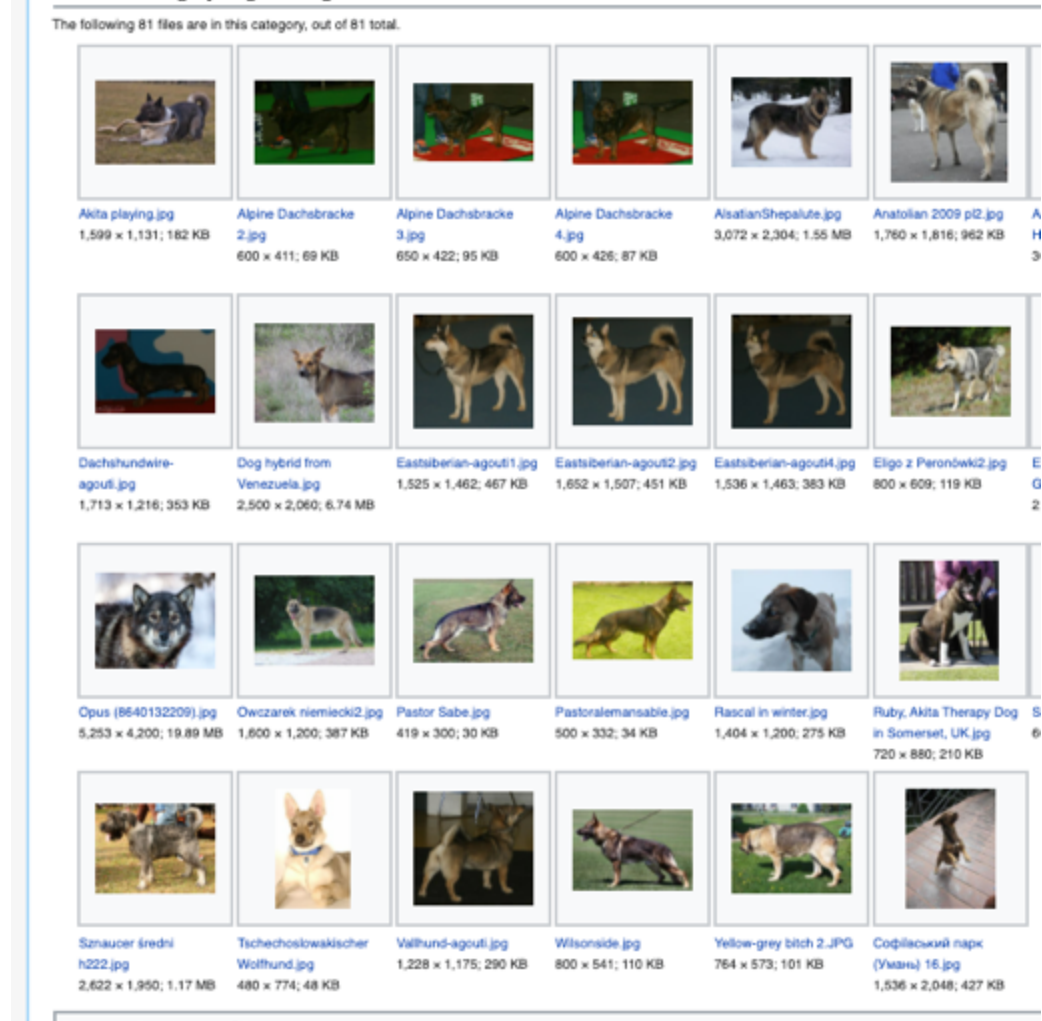
| Akita playing.jpg | Alpine Dachsbracke 2.jpg | Alpine Dachsbracke 3.jpg | Alpine Dachsbracke 4.jpg | AlsatianShepalute.jpg | Anatolian 2009 pl2.jpg |
|---|---|---|---|---|---|
| 1,599 × 1,131; 182 KB | 600 × 411; 69 KB | 650 × 422; 95 KB | 600 × 426; 87 KB | 3,072 × 2,304; 1.55 MB | 1,760 × 1,816; 962 KB |

| Dachshundwire-agouti.jpg | Dog hybrid from Venezuela.jpg | Eastsiberian-agouti1.jpg | Eastsiberian-agouti2.jpg | Eastsiberian-agouti4.jpg | Eligo z Peronówki2.jpg |
|---|---|---|---|---|---|
| 1,713 × 1,216; 353 KB | 2,500 × 2,060; 6.74 MB | 1,525 × 1,462; 467 KB | 1,652 × 1,507; 451 KB | 1,536 × 1,463; 383 KB | 800 × 609; 119 KB |

| Opus (8640132209).jpg | Owczarek niemiecki2.jpg | Pastor Sabe.jpg | Pastoralemansable.jpg | Rascal in winter.jpg | Ruby, Akita Therapy Dog in Somerset, UK.jpg |
|---|---|---|---|---|---|
| 5,253 × 4,200; 19.89 MB | 1,600 × 1,200; 387 KB | 419 × 300; 30 KB | 500 × 332; 34 KB | 1,404 × 1,200; 275 KB | 720 × 880; 210 KB |

| Sznaucer średni h222.jpg | Tschechoslowakischer Wolfhund.jpg | Vallhund-agouti.jpg | Wilsonside.jpg | Yellow-grey bitch 2.JPG | Софіївський парк (Умань) 16.jpg |
|---|---|---|---|---|---|
| 2,622 × 1,950; 1.17 MB | 480 × 774; 48 KB | 1,228 × 1,175; 290 KB | 800 × 541; 110 KB | 764 × 573; 101 KB | 1,536 × 2,048; 427 KB |

$\mathcal{D}^{(1)}$
learn h by fine-tuning $\hat{h}$

$\mathcal{D}^{(2)}$
pre-train hypothesis $\hat{h}$

$$\min_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}\big((x^{(i)}, y^{(i)}), h\big) + \lambda d(h, \hat{h})$$
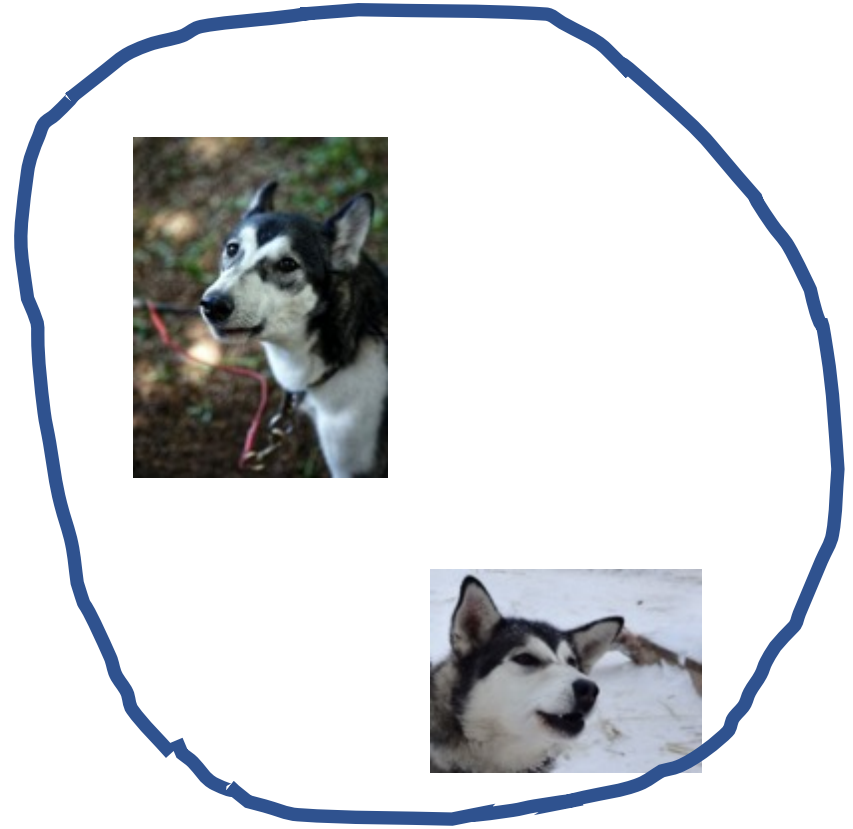
fine tuning on $\mathcal{D}^{(1)}$

distance to hypothesis $\hat{h}$ which is pre-trained on $\mathcal{D}^{(2)}$

# Multi-Task Learning via Regularization

- Problem I: classify image as "shows border colly" vs. "not"

- Problem II: classify image as "shows husky" vs. "not"

- training data $\mathcal{D}^{(1)}$ for Problem I and $\mathcal{D}^{(2)}$ for Problem II

- jointly learn hypothesis $h^{(1)}$ on $\mathcal{D}^{(1)}$ and $h^{(2)}$ on $\mathcal{D}^{(2)}$

- require $h^{(1)}$ to be "similar" to $h^{(2)}$

$$\mathcal{D}^{(1)} \qquad\qquad \mathcal{D}^{(2)}$$

jointly learn similar
$h^{(1)}$ and $h^{(2)}$ for each dataset

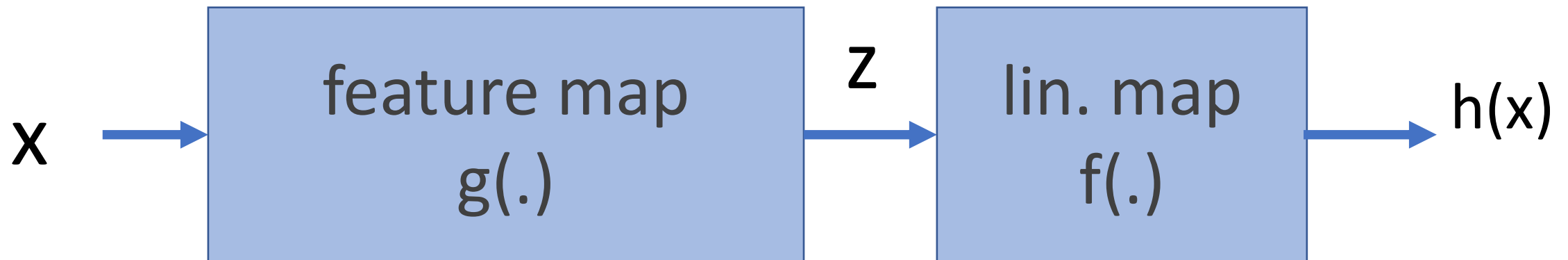training error of $h^{(1)}$

training error of $h^{(2)}$

$$\min_{h^{(1)}, h^{(2)}} \quad \mathcal{E}\big(h^{(1)}\big|\mathcal{D}^{(1)}\big) + \mathcal{E}\big(h^{(2)}\big|\mathcal{D}^{(2)}\big) + \lambda d\big(h^{(1)}, h^{(2)}\big)$$
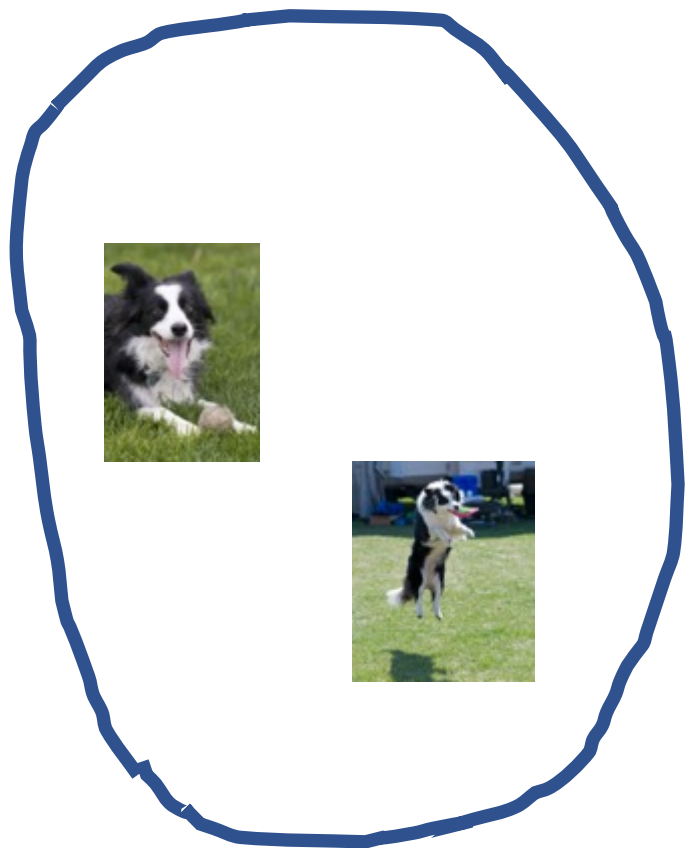
"distance" between $h^{(1)}$ and $h^{(2)}$

# Semi-Supervised Learning via Regularization

- classify image as "shows border colly" vs. "not"

- small labeled dataset $\mathcal{D}^{(1)}$

- massive image database $\mathcal{D}^{(2)}$ with unlabeled images

- train hypothesis h(.) on $\mathcal{D}^{(1)}$ with following structure:

X → feature map g(.) → z → lin. map f(.) → h(x)

$$\mathcal{D}^{(1)}$$

learn linear classifier f(.)

$$\mathcal{D}^{(2)}$$

learn feature map g(.)

$$\min_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}\left((x^{(i)}, y^{(i)}), h\right) + \lambda \, \mathcal{E}\left(g | \mathcal{D}^{(2)}\right)$$

use training error
to fine tune f(.)

learn feature map g(.)
using large unlabeled
database $\mathcal{D}^{(2)}$

# To sum up,

- regularization is a soft model pruning

- regularization does implicit data augmentation

- special cases of regularization

  - transfer learning

  - multi-task learning

  - semi-supervised learning

# Questions ?