# Regularization

Alex(ander) Jung
Assistant Professor for Machine Learning
Department of Computer Science
Aalto University

# Reading.

Ch. 7 of https://mlbook.cs.aalto.fi

# Learning Goals

- develop intuition for <span style="color:red">effective data and model size</span>

- <span style="color:red">reduce model size</span> by <span style="color:red">model pruning</span>

- <span style="color:red">increase data size</span> by <span style="color:red">data augmentation</span>

- regularization = impl. model pruning = impl. data aug.

- use reg. for transfer - , multi-task – and semi-supervised learning

# Empirical Risk Minimization

learn hypothesis out of model that incurs minimum loss when predicting labels of datapoints based on their features

training set

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

(2.16)
$$= \operatorname*{argmin}_{h \in \mathcal{H}} (1/m) \sum_{i=1}^{m} L\big((\mathbf{x}^{(i)}, y^{(i)}), h\big).$$
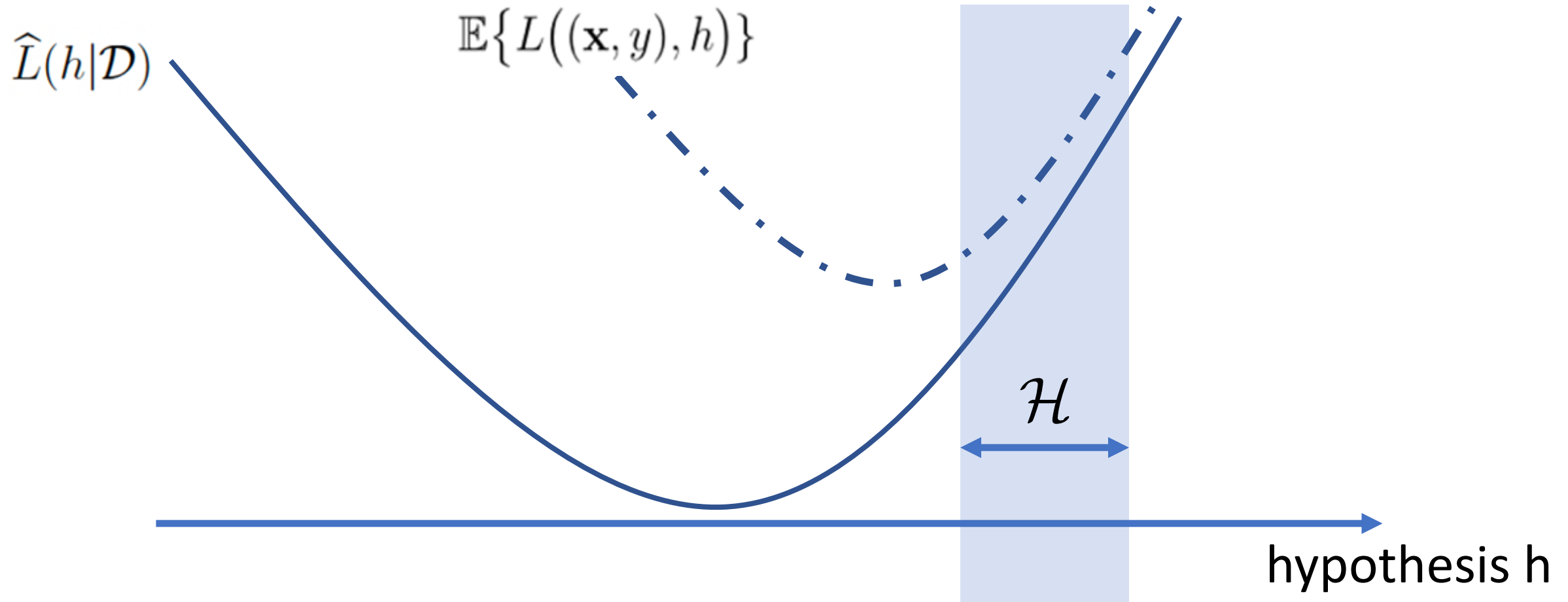
hypothesis

label of i-th datapoint

features of i-th datapoint

model
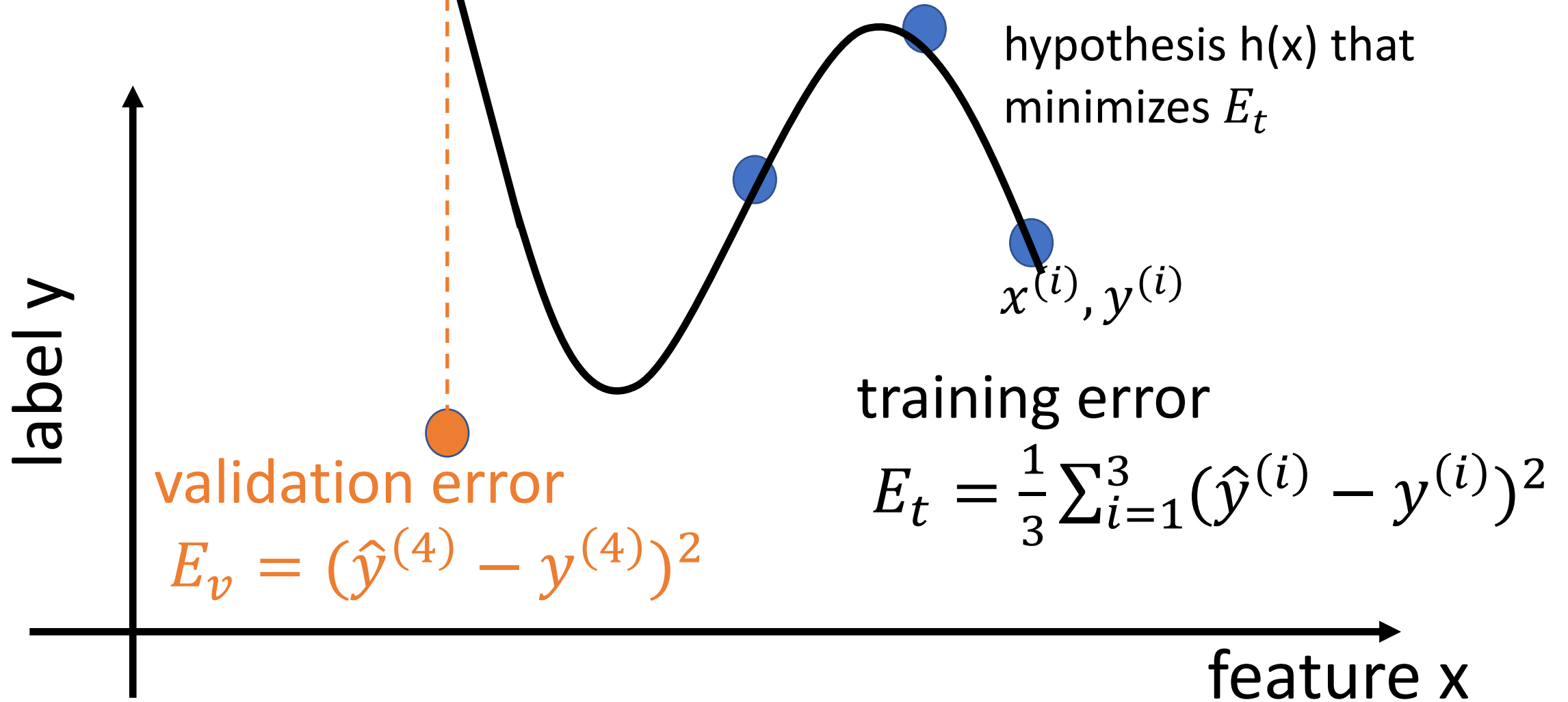
loss function

see Ch. 4.1 of mlbook.cs.aalto.fi

# ERM is only Approximation!

$\widehat{L}(h|\mathcal{D})$

$\mathbb{E}\{L((\mathbf{x}, y), h)\}$

$\mathcal{H}$

hypothesis h

# Train and Validate Model $\mathcal{H}^{(3)}$



hypothesis h(x) that minimizes $E_t$

$x^{(i)}, y^{(i)}$

training error

$$E_t = \frac{1}{3}\sum_{i=1}^{3}(\hat{y}^{(i)} - y^{(i)})^2$$

validation error

$$E_v = (\hat{y}^{(4)} - y^{(4)})^2$$

label y

feature x

# Small Training Error Does Not Imply Good Performance on New Data Points!

# One Pixel Attack for Fooling Deep Neural Networks
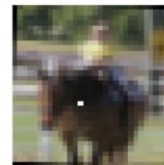
Jiawei Su*, Danilo Vasconcellos Vargas* and Kouichi Sakurai

rch has revealed that the output of Deep
an be easily altered by adding relatively
input vector. In this paper, we analyze
limited scenario where only one pixel
t we propose a novel method for gen-
ial perturbations based on differential
s less adversarial information (a black-
l more types of networks due to the
The results show that 67.97% of the
e CIFAR-10 test dataset and 16.04%
C 2012) test images can be perturbed
ass by modifying just one pixel with
idence on average. We also show the
original CIFAR-10 dataset. Thus, the
a different take on adversarial machine
imited scenario, showing that current
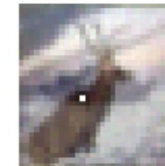


**AllConv**
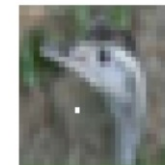SHIP
CAR(99.7%)

HORSE
DOG(70.7%)

**NiN**
HORSE
FROG(99.9%)

DOG
CAT(75.5%)

**VGG**
DEER
AIRPLANE(85.3%)
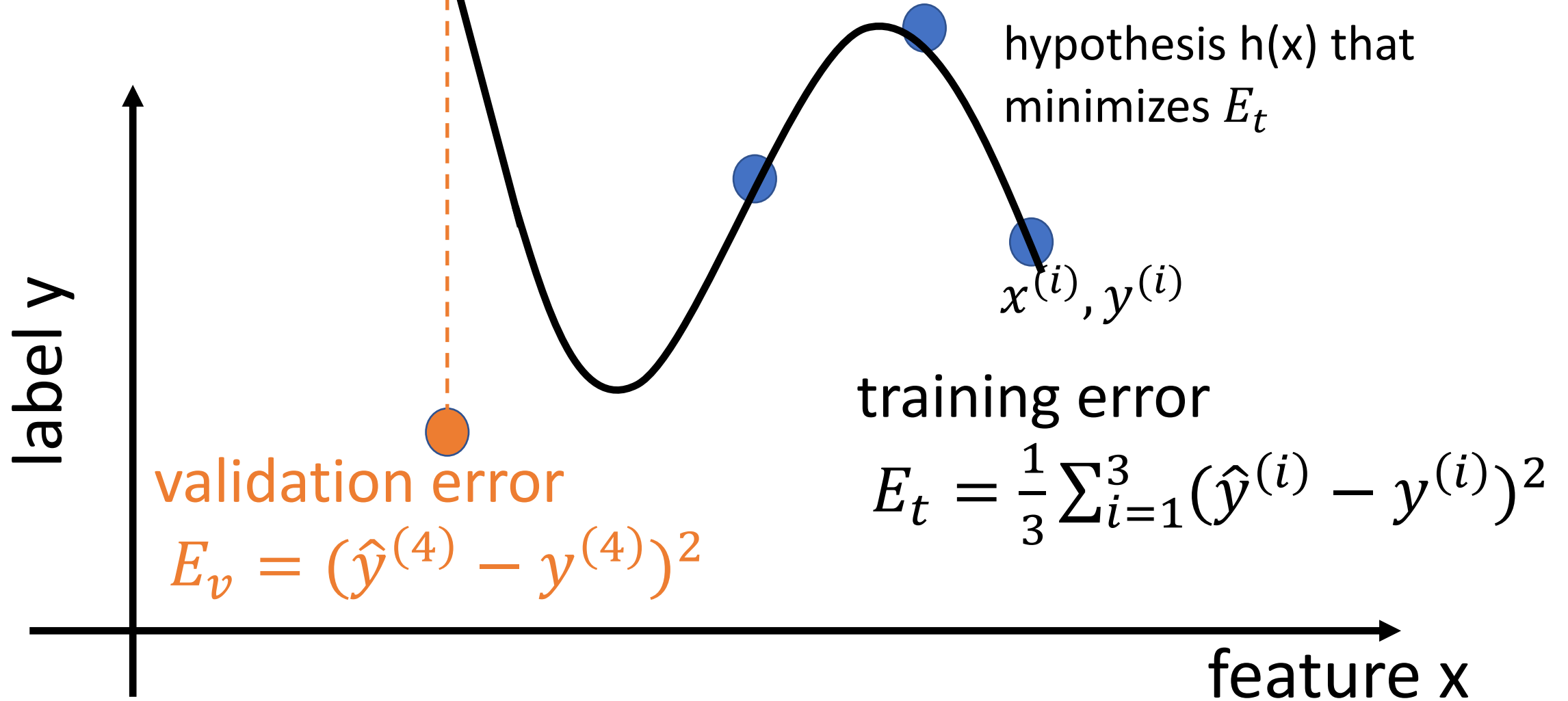
BIRD
FROG(86.5%)

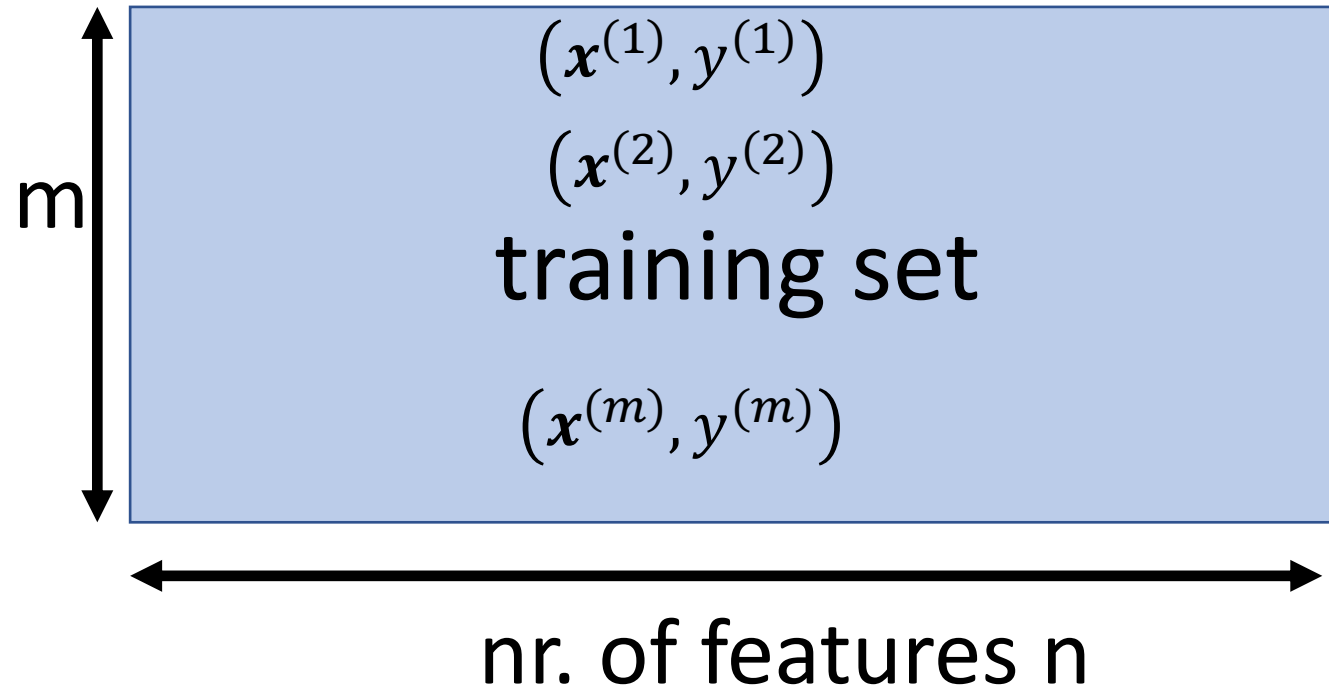https://arxiv.org/pdf/1710.08864.pdf

# EU Guidelines for Trustworthy AI

"…Technical Robustness and safety: AI systems need to be resilient and secure. They need to be safe, ensuring a fall back plan in case something goes wrong, as well as being accurate, reliable and reproducible. That is the only way to ensure that also unintentional harm can be minimized and prevented…."

https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai

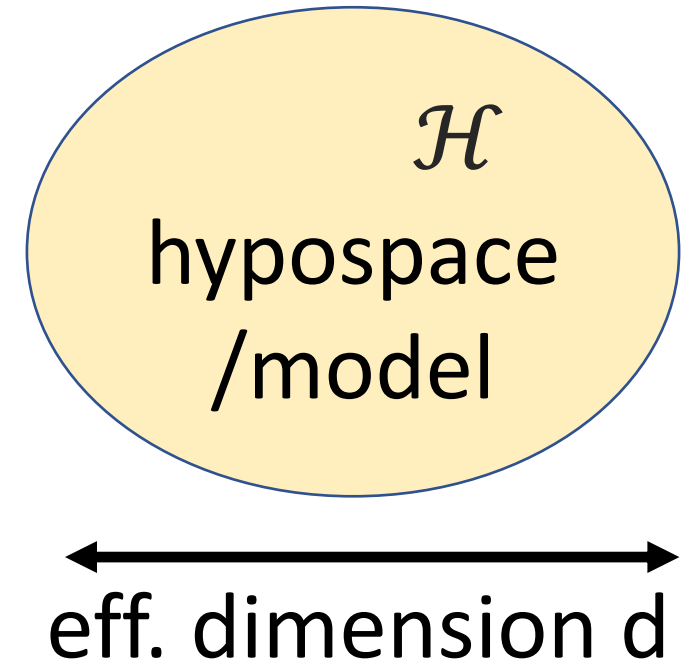# Single Pixel Attacks !

hypothesis h(x) that minimizes $E_t$

$x^{(i)}, y^{(i)}$

training error

$$E_t = \frac{1}{3}\sum_{i=1}^{3}(\hat{y}^{(i)} - y^{(i)})^2$$

validation error

$$E_v = (\hat{y}^{(4)} - y^{(4)})^2$$

label y

feature x

# Data and Model Size

$(\boldsymbol{x}^{(1)}, y^{(1)})$

$(\boldsymbol{x}^{(2)}, y^{(2)})$

m

training set

$(\boldsymbol{x}^{(m)}, y^{(m)})$

nr. of features n

$\mathcal{H}$
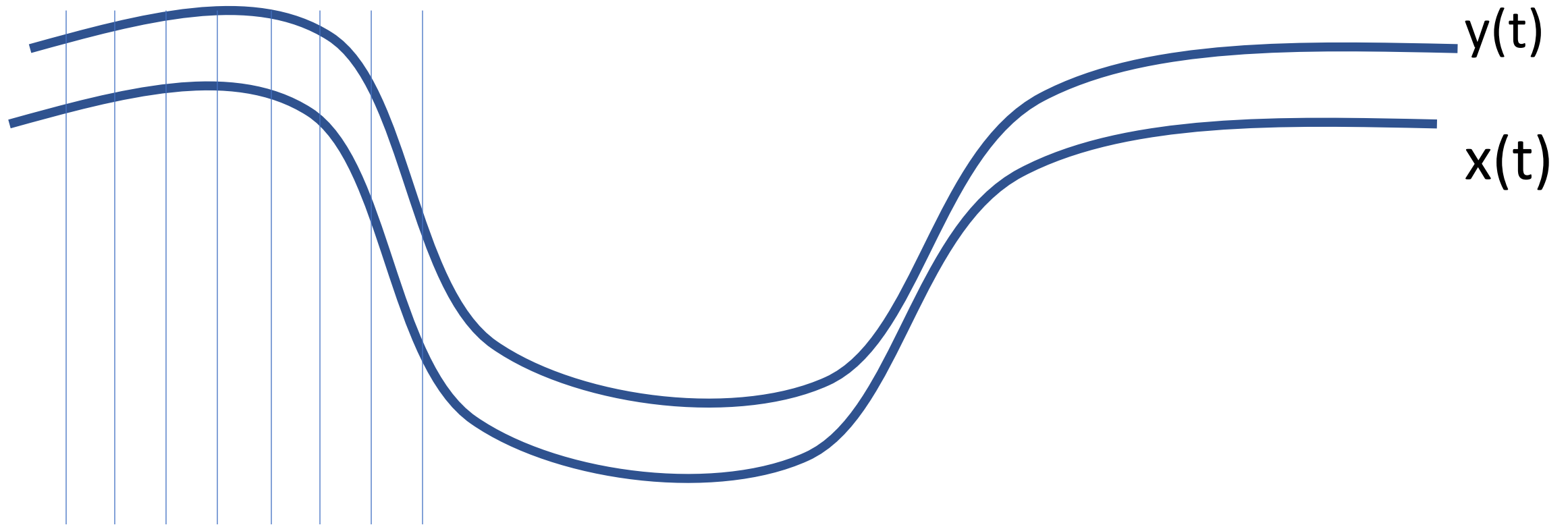hypospace /model

eff. dimension d

crucial parameter is the ratio d/m

# Effective Data Size

consider data points obtained from time series
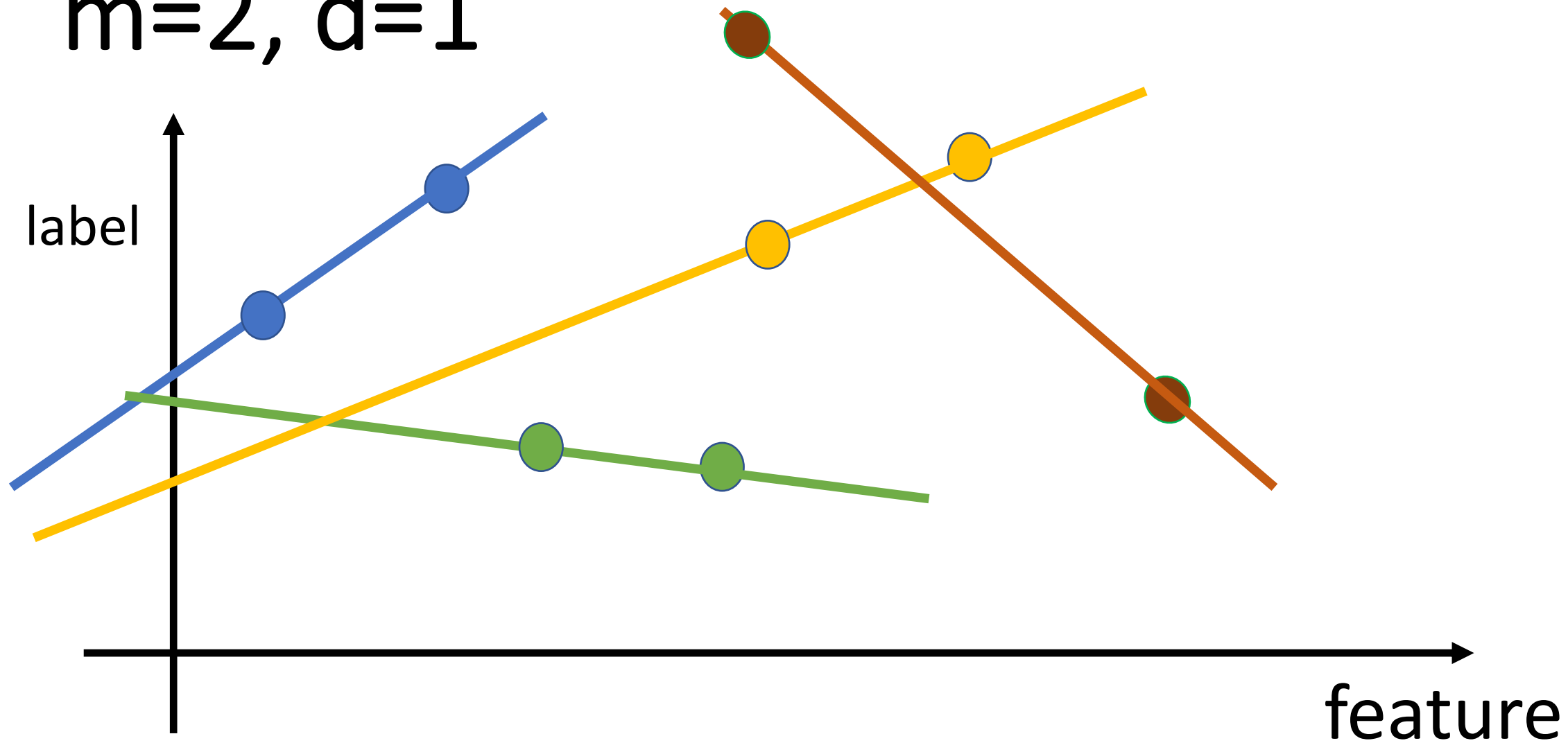


y(t)

x(t)

# Effective Dim. Linear Maps

- linear map can perfectly fit m data points with n features, as soon as n ≥ m [Ch 6.1, mlbook.cs.aalto.fi]

- eff.dim. of linear maps = nr. of features
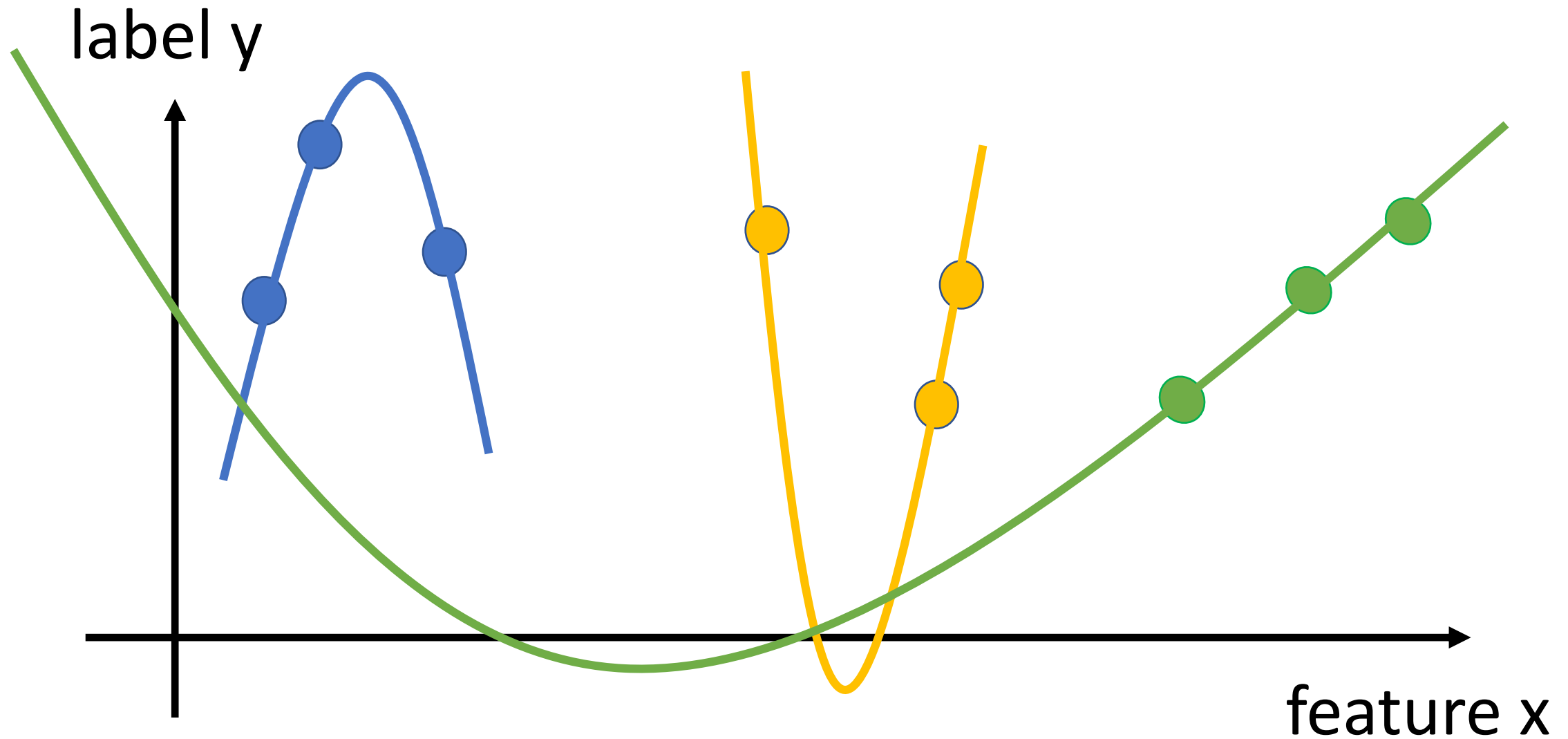
- d = n

# Effective Dim. Linear Maps

we can perfectly fit (almost) any <span style="color:red">m data points</span> using polynomials of <span style="color:red">degree d</span> as soon as

$$d \geq m-1$$

# m=2, d=1

label

feature

# m=3, degree d=2 polynomial

# Data Hungry ML Methods

- millions of features for datapoints (e.g. megapixel image)

- eff.dim. d of linear maps is also millions

- eff.dim d of deep nets is millions ... billions

- can perfectly fit any set of 100000s (!) of datapoints

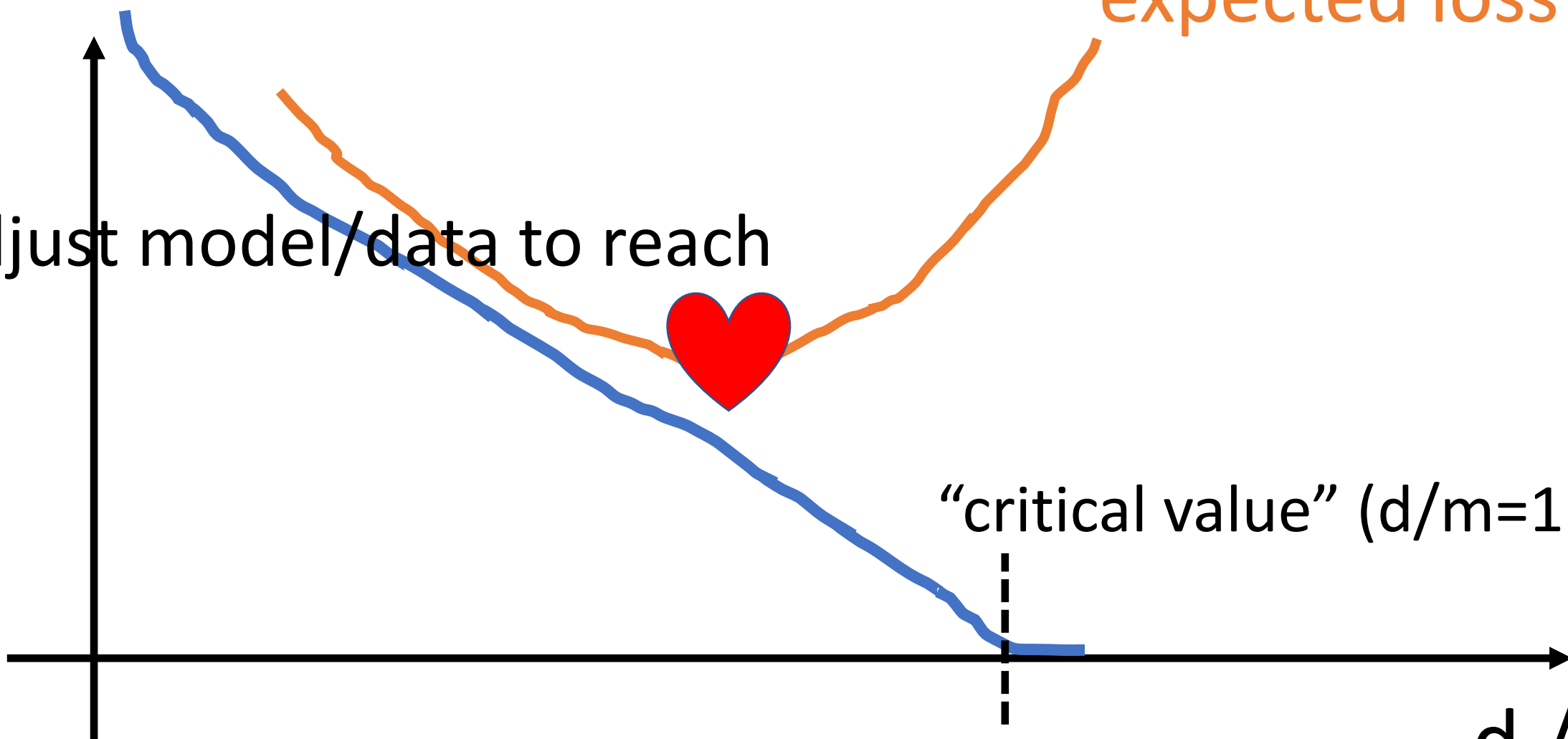- training error will be zero (overfitting!)

training error

expected loss

adjust model/data to reach

"critical value" (d/m=1)

d / m

# how to bring d/m below critical value?
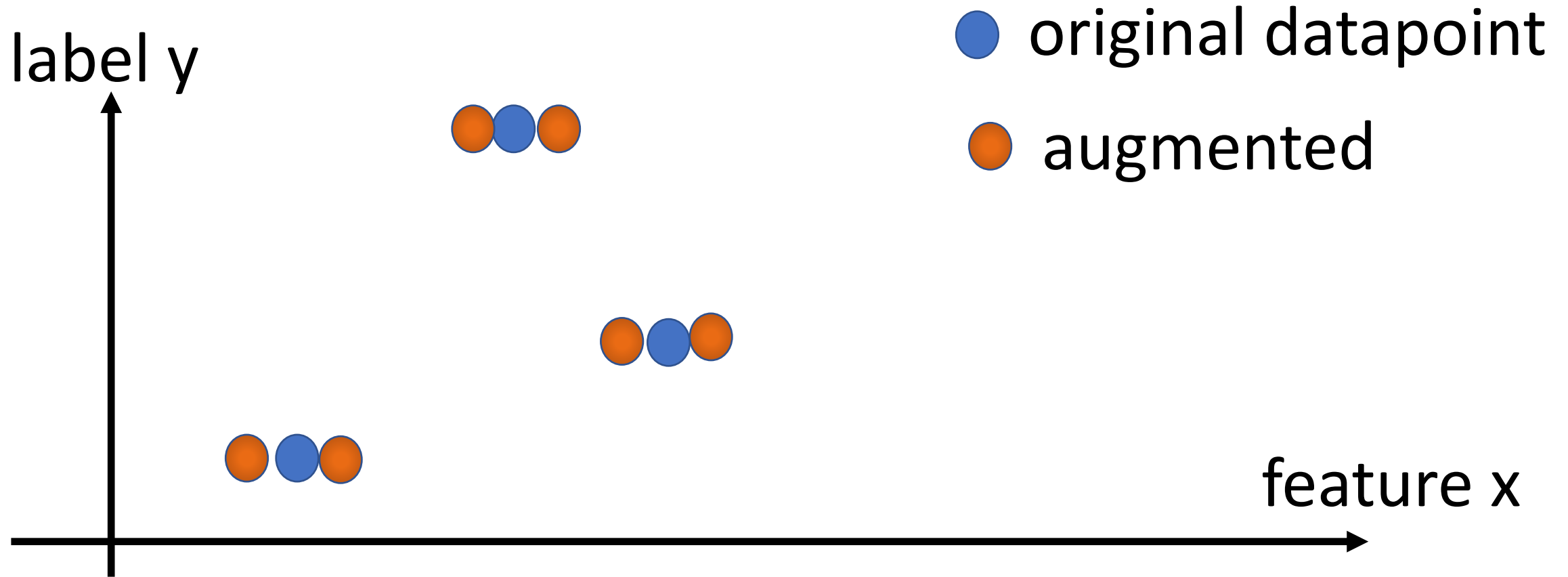
- increase m by using more training data

- decrease d by using smaller hypothesis space

# how to bring d/m below critical value?

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

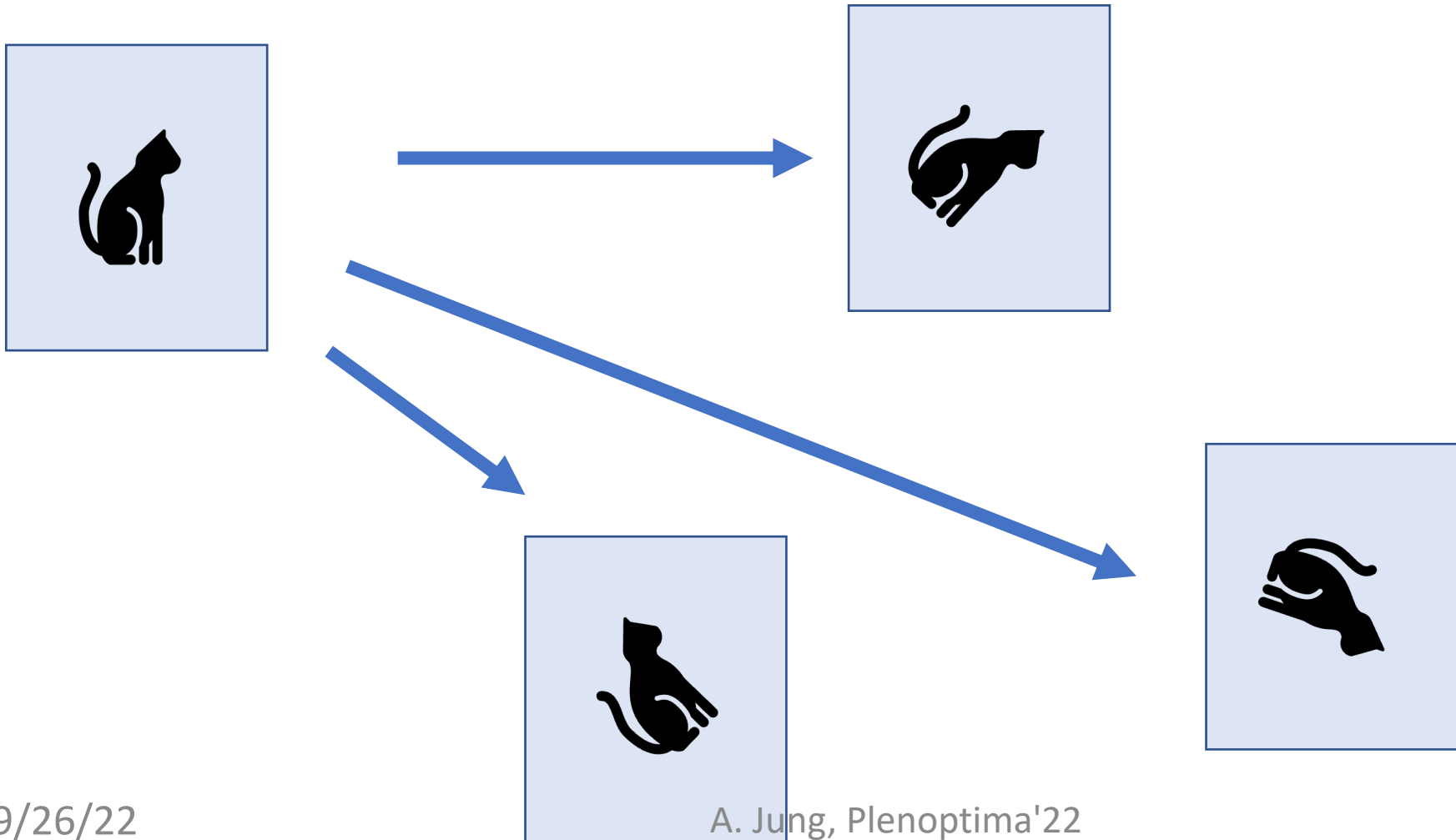- decrease d by using smaller hypothesis space

# 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

A. Jung, Plenoptima'22

# flipped cat image is still cat image



A. Jung, Plenoptima'22

# shifted cat image is still cat image



A. Jung, Plenoptima'22

# noisy cat image is still cat image

A. Jung, Plenoptima'22

# how to bring d/m below critical value?

- increase m by using more training data

- <span style="color:red">decrease d by using smaller hypothesis space</span>
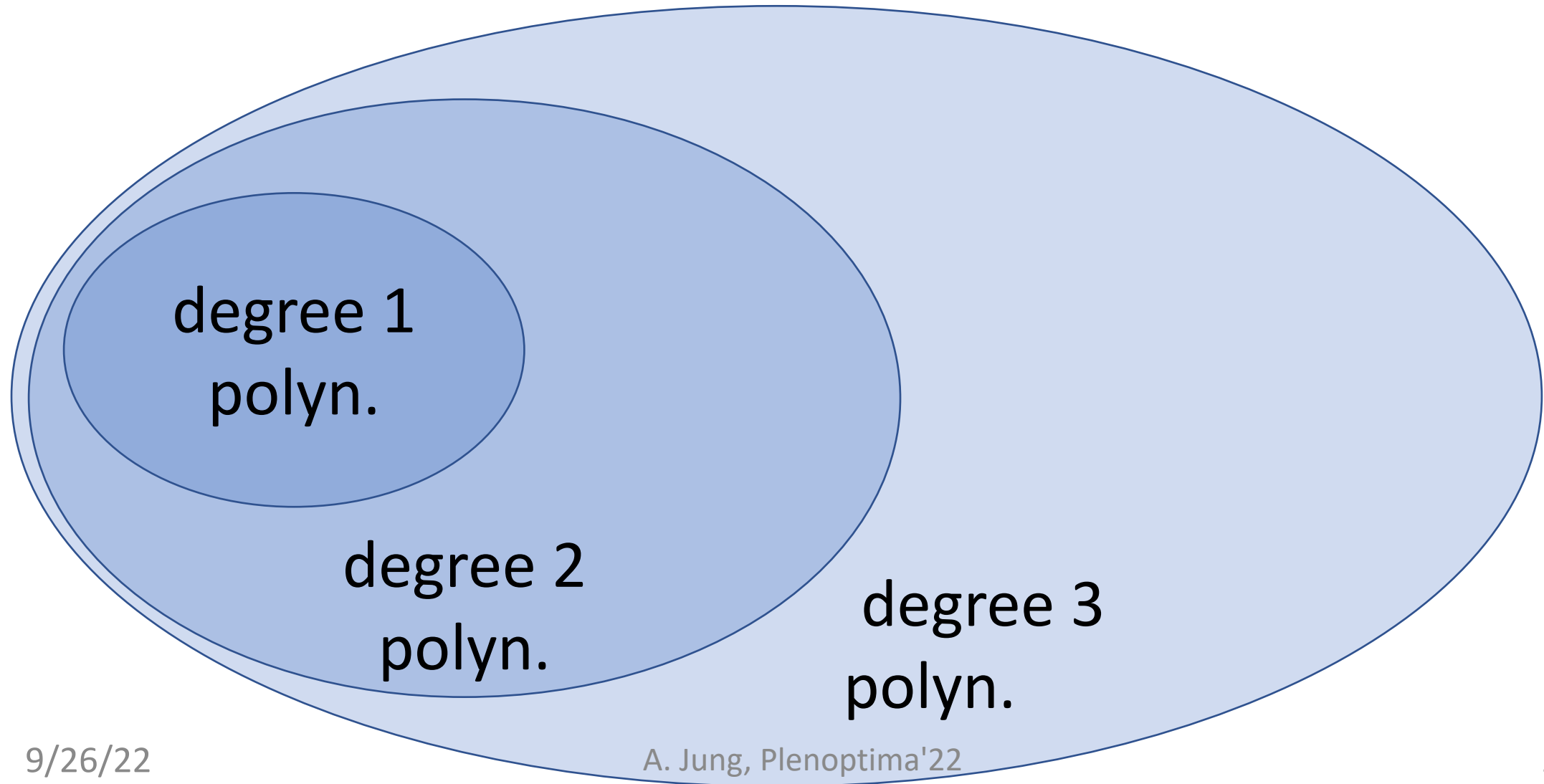
replace original ERM

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

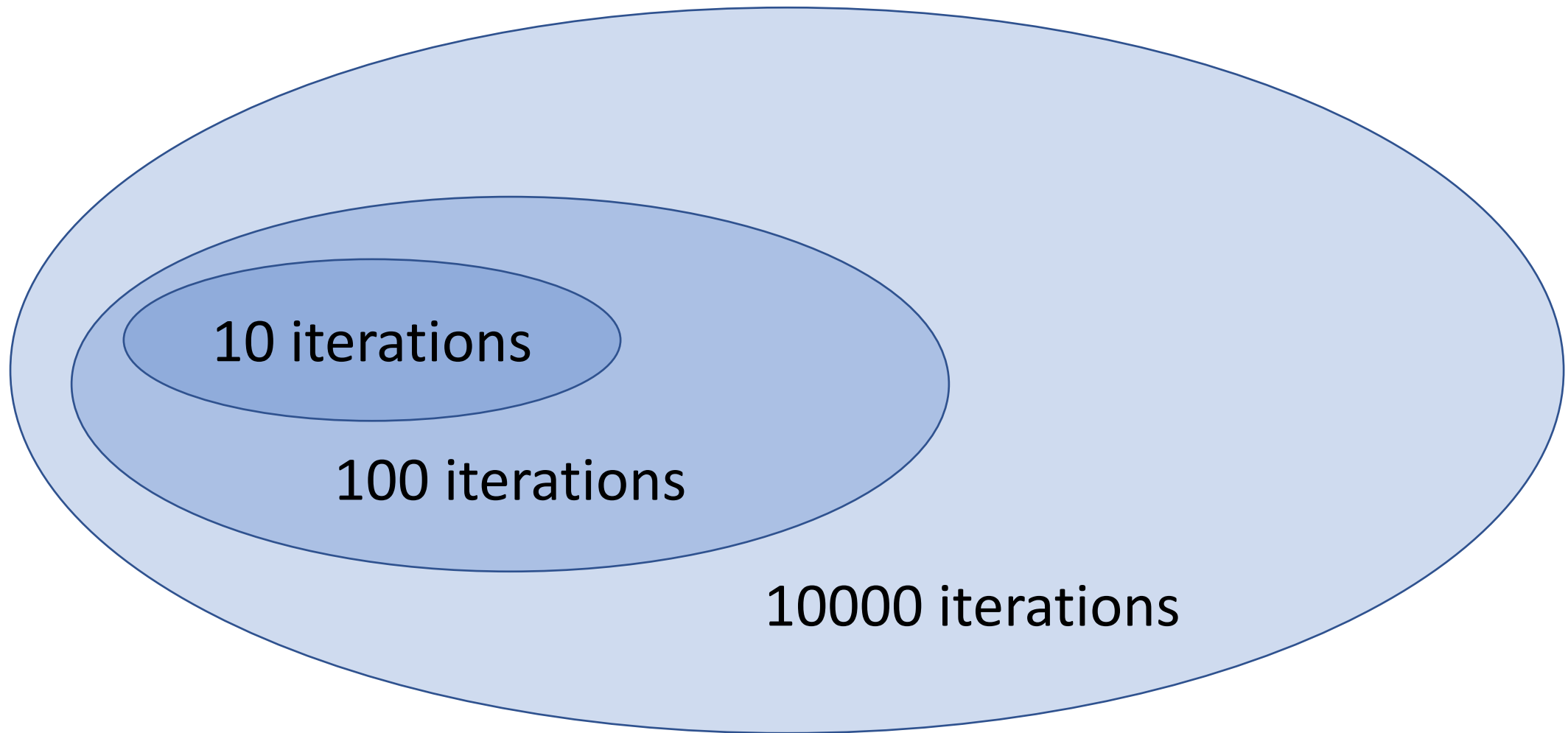with ERM on smaller $\widehat{\mathcal{H}} \subset \mathcal{H}$

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

# Nested Models



degree 1 polyn.

degree 2 polyn.

degree 3 polyn.

A. Jung, Plenoptima'22

# Prune Hypospace by Early Stopping



10 iterations

100 iterations

10000 iterations

# Soft Model Pruning via Regularization

A. Jung, Plenoptima'22

# Regularized ERM

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

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

average loss on training set
(empirical risk of h)

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 $\mathcal{R}(w) = \|w\|_1 = |w_1| + \cdots + |w_n|$
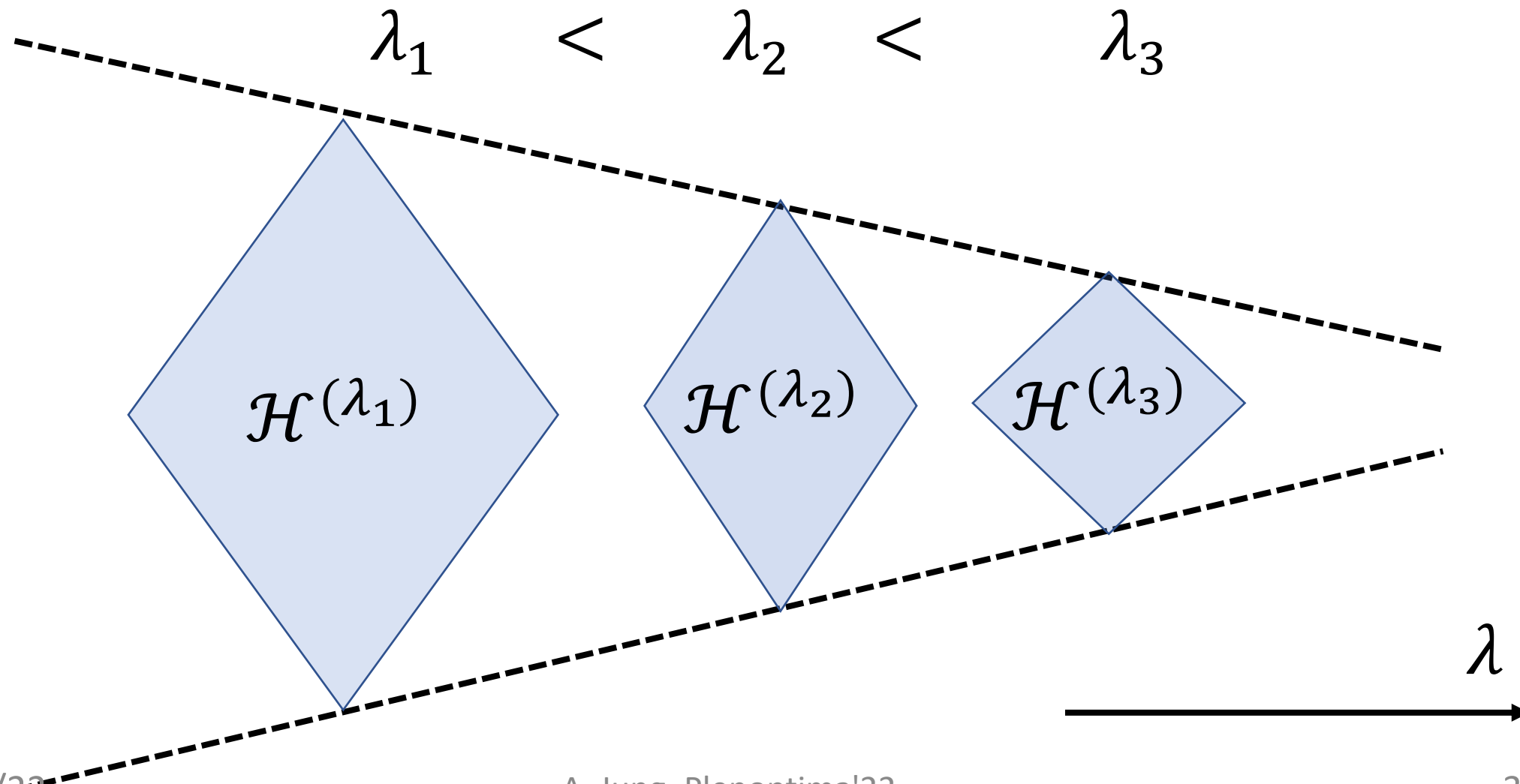
# Regularization = Implicit Pruning!

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

## equivalent to

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

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

A. Jung, Plenoptima'22

# Regularization = "Soft" Model Selection

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



$$\mathcal{H}^{(\lambda_1)} \qquad \mathcal{H}^{(\lambda_2)} \qquad \mathcal{H}^{(\lambda_3)}$$

$\lambda$

# Regularization does implicit Data Augmentation

# augment with (infinitely many) realizations of RV!

label y

● original datapoint

● augmented

● = ● + "noise"

feature x

A. Jung, Plenoptima'22

# Regularization =Implicit Data Aug.

label y

raw datapoint

"perturbed"
datapoint

h(x)

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

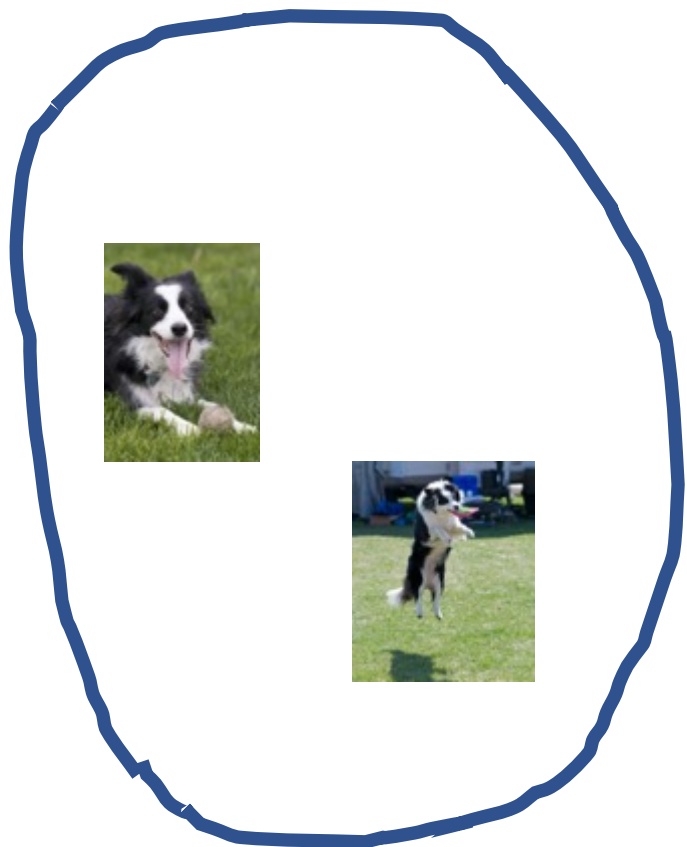see Chapter 7.3 of mlbook.cs.aalto.fi

feature x

# To sum up,

- large ratio d/m leads to overfitting

- reduce d by using smaller model ("pruning")

- increase m by using more data points

- regularization is a soft model pruning

- regularization does implicit data augmentation

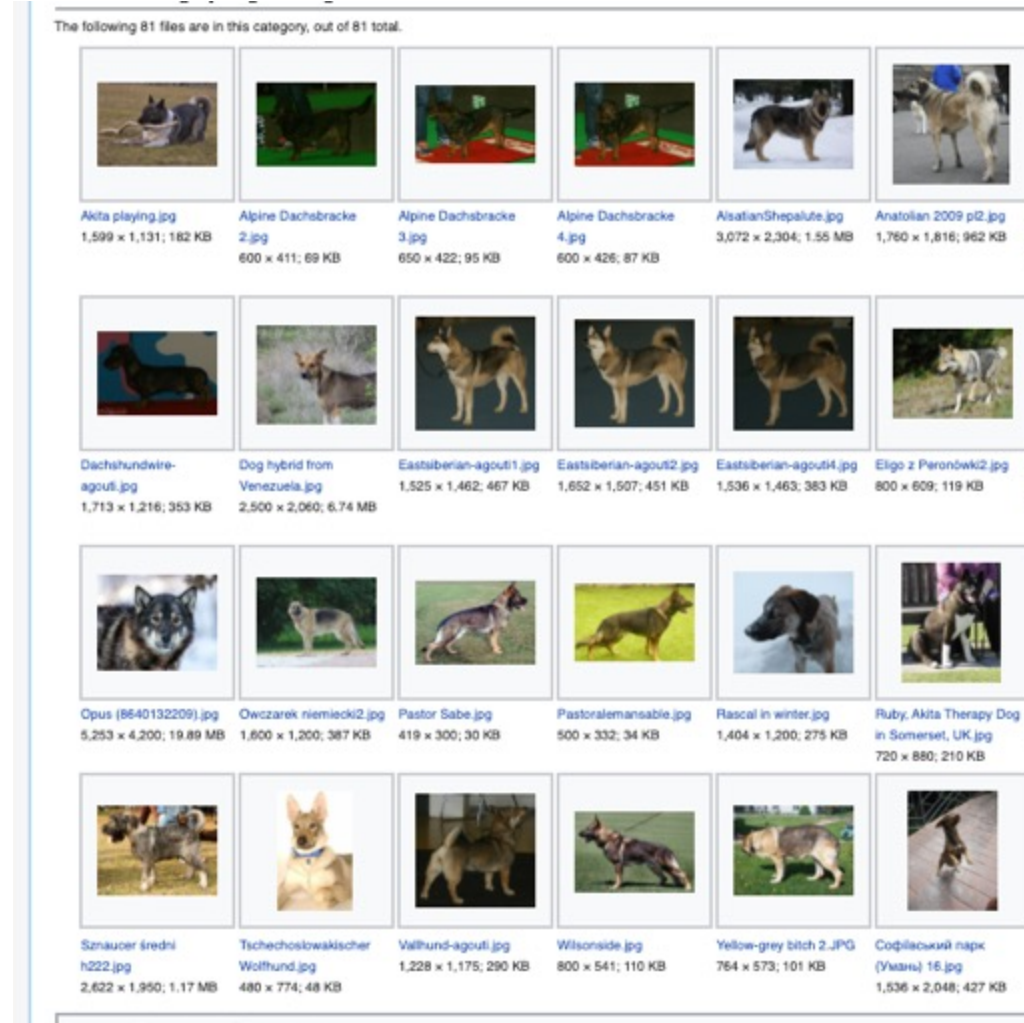# Transfer Learning via Regularization

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

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

learn h by fine-tuning $\hat{h}$

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

pre-train hypothesis $\hat{h}$

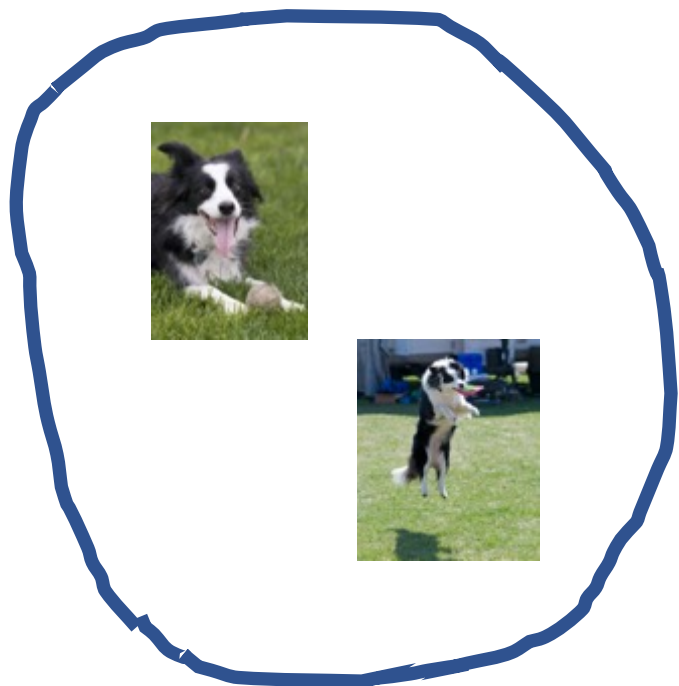$$\min_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^{m} L\big((x^{(i)}, y^{(i)}), h\big) + \lambda d(h, \hat{h})$$
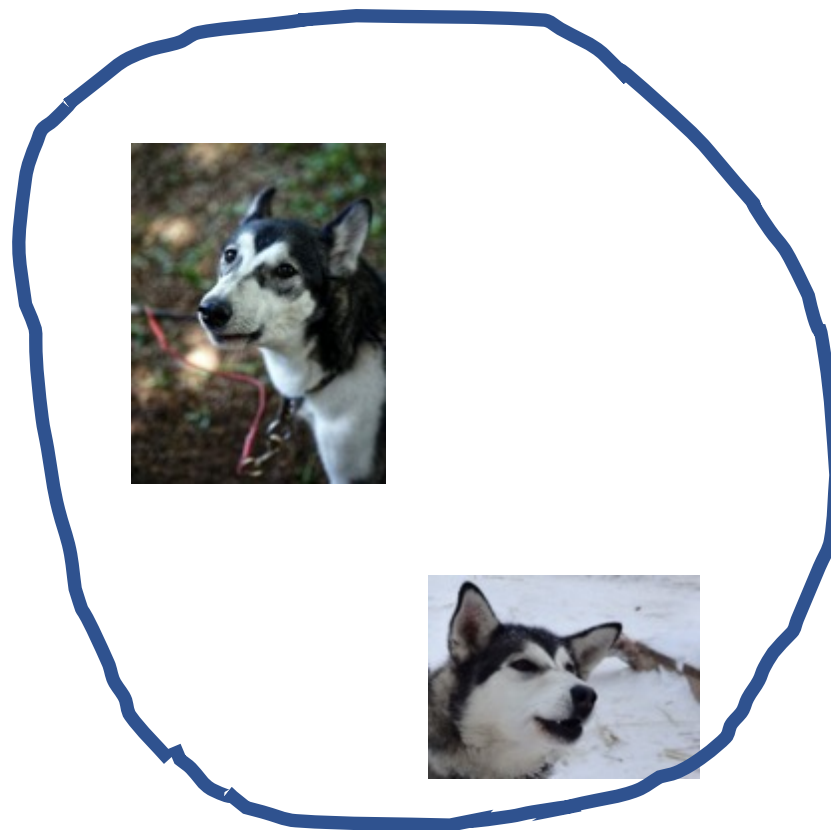
fine tuning on $\mathcal{D}^{(1)}$

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

# Multi-Task Learning via Regularization

A. Jung, Plenoptima'22

- 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

- <span style="color:red">jointly learn</span> hypothesis $h^{(1)}$ on $\mathcal{D}^{(1)}$ and $h^{(2)}$ on $\mathcal{D}^{(2)}$

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

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

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

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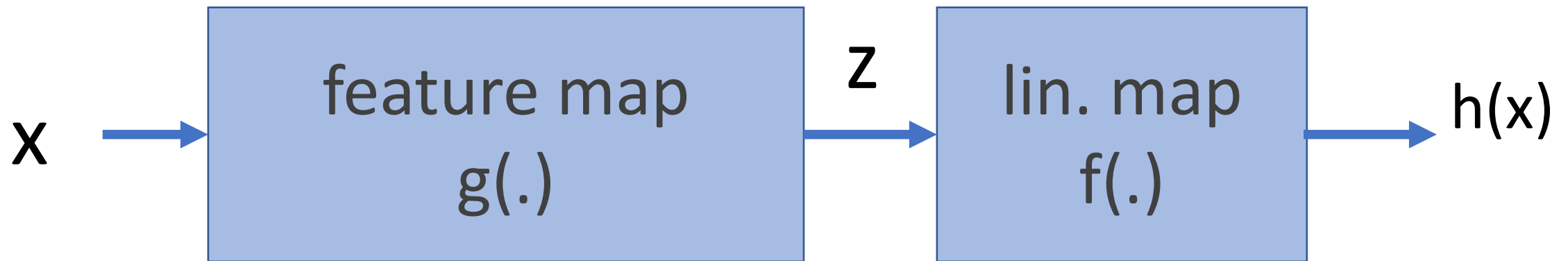
training error of $h^{(1)}$

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

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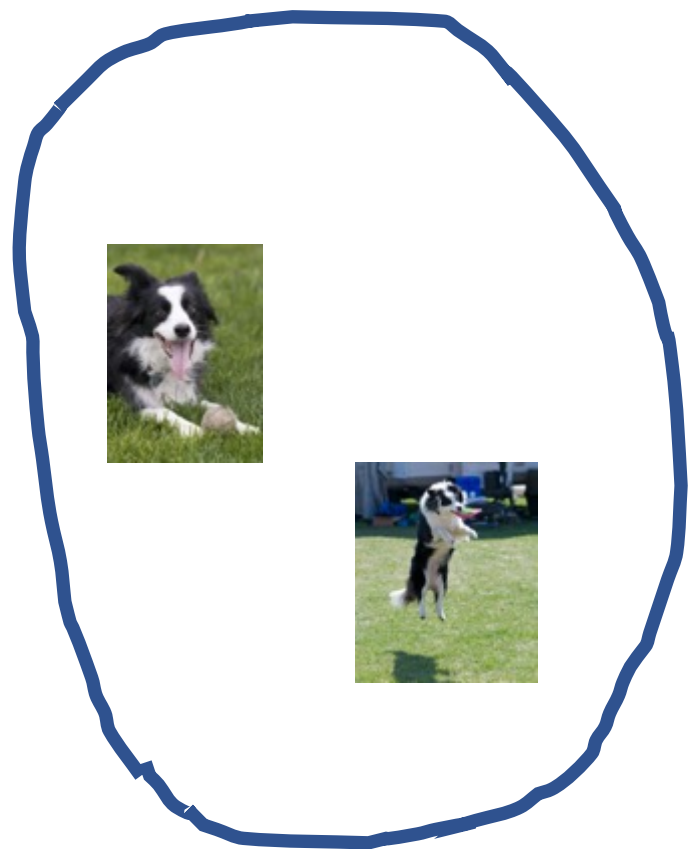
"distance" between $h^{(1)}$ and $h^{(2)}$

A. Jung, Plenoptima'22

# 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)

"chain" or "pipeline"

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

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

learn linear classifier f(.)    learn feature map g(.)

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

use training error
to fine tune h(.)

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

A. Jung, Plenoptima'22

# Subjective Explainability via Regularization

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

- $\mathrm{E}(h|u)$ measures explainability of hypothesis h(.) to user u

- want same h(x) for data points with similar user signal u

- implementation of "Human agency and oversight"

# To Sum Up

- ML works well if m/d > 1

- increase data size m by data augmentation

- decrease model size d by regularization

- adding reg. term = data augmentation/soft model-pruning

- transfer-, multi-task- and semi-supervised learning as instances of regularization