# Best Practice in Machine Learning for Computer Vision

Christoph H. Lampert

Institute of Science and Technology Austria (ISTA)

Vision and Sports Summer School 2022

July 25-29, 2022 — Prague, Czech Republic

**I S T AUSTRIA**

*Institute of Science and Technology*

Slides and exercise data: https://cvml.ist.ac.at/

**Research topics in ISTA's ELLIS Unit:**
- ▶ transfer/trustworthy learning, computer vision (me)
- ▶ scalable deep learning (Alistarh group)
- ▶ theory of deep learning (Mondelli group)
- ▶ Bayesian methods in genomics (Robinson group)
- ▶ machine learning for material science (Cheng group)
- ▶ discrete optimization (Kolmogorov group)

**ISTA Graduate School**
- ▶ two-phase,English language PhD program
- ▶ full salary positions

**Other possibilities: internships, postdocs, faculty, . . .**
- ▶ visit: `cvml.ist.ac.at` or ask me during a break

**Computer Vision**

**Machine Learning for Computer Vision**

**Some Background**

**Best Practice and How to Avoid Pitfalls**

# Computer Vision

# CIFAR-10 - Object Recognition in Images

Identify the subject of 60,000 labeled images

Kaggle · 231 teams · 8 years ago

Overview    Data    Code    Discussion    Leaderboard    Rules
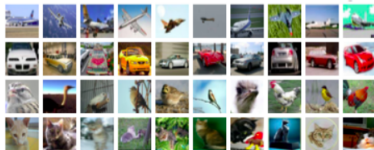
**Join Competition**    ...

## Overview

Description

Evaluation

CIFAR-10 is an established computer-vision dataset used for object recognition. It is a subset of the 80 million tiny images dataset and consists of 60,000 32×32 color images containing one of 10 object classes, with 6000 images per class. It was collected by Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton.

Kaggle is hosting a CIFAR-10 leaderboard for the machine learning community to use for fun and practice. You can see how your approach compares to the latest research methods on Rodrigo Benenson's classification results page.

nd Prediction Competition

# CI... ...ject Recognition in Images

Identify the ... ...0 labeled images

Kaggle · 231 teams ·

Overview  Data  Code  Discussion  Le...

**Join Competition**  ...

## Overview

**Description**

Evaluation

CIFAR-10 is an esta... ...tion dataset ... ...tition. It is a subset of the 80 million tiny images dataset and co... ...2 color images cont... ...lasses, with 6000 images per class. It was collected ... ...nod Nair, and Geoffrey Hinton.

K... ...AR-10 leaderboard for the machine learning com... ...d practice. You can see how your ... ...es to the latest research methods on Rodrigo Benenson's cla...

Robotics (e.g. *autonomous cars*)



Healthcare (e.g. *visual aids*)



Commerce (e.g. *Amazon Go)*)



Consumer Electronics
(e.g. *human computer interaction*)



Augmented Reality
(e.g. *HoloLens, Pokemon Go*)



Security
(e.g. *Face Unlock*)

Amazon Go Store image: amazon.com

Computer vision systems should

- ▶ perform interesting/relevant tasks,
  e.g. drive a car

For this, they need to

- ▶ work in the real world,
- ▶ interact with non-experts,
- ▶ do what they are supposed to do
  without failures.
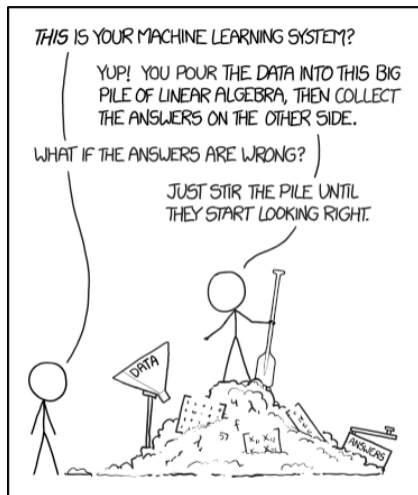
Computer vision systems should

▶ perform interesting/relevant tasks,
   e.g. drive a car

For this, they need to

▶ work in the real world,
▶ interact with non-experts,
▶ do what they are supposed to do
   without failures.

**Machine learning is not particularly
good at those...**

▶ we have to be extra careful to know
   what we're going!



[Image: https://xkcd.com/1838/]

# Machine Learning for Computer Vision

## Example: Solving a Computer Vision Task using Machine Learning

Step 0) Sanity checks...

Step 1) Decide what exactly you want
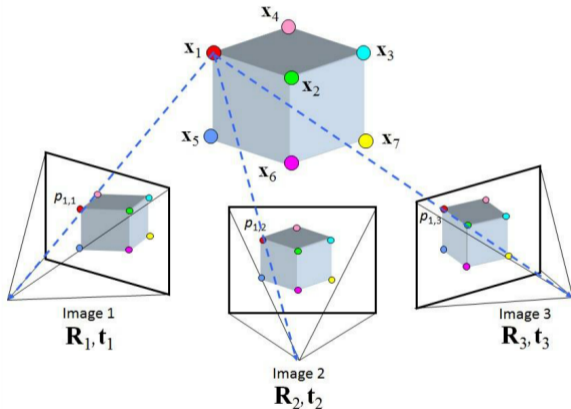
Step 2) Collect and annotate data

Step 3) Model training

Step 4) Model evaluation

Step 0) Sanity checks...

Step 1) Decide what exactly you want

Step 2) Collect and annotate data

Step 3) Model training

Step 4) Model evaluation

This lecture:

▶ **Question 1: why do we do it like this?**

▶ **Question 2: what can go wrong, and how to avoid that?**

**Create a 3D model from many 2D correspondences**



Sanity Check: Is Machine Learning the right way to solve it?

▶ Can you think of an algorithmic way to solve the problem?

   **Yes:** optimization with geometric constraints!

$\rightarrow$ not a strong case for using machine learning. ✗

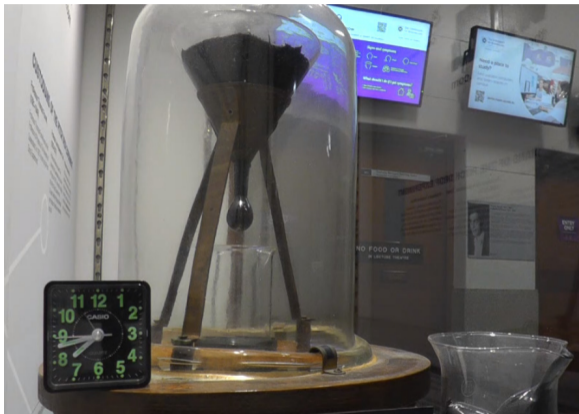images: http://vision.in.tum.de/research/image-based_3d_reconstruction/multiviewreconstruction

Sanity Check: Is Machine Learning the right way to solve it?

► Can you think of an algorithmic way to solve the problem?    **No.**

image: http://www.thetenthwatch.com/feed/

Sanity Check: Is Machine Learning the right way to solve it?

► Can you think of an algorithmic way to solve the problem?  **No.**

► Can you get sufficient data? **No.**

$\rightarrow$ not a strong case for using machine learning. ✗
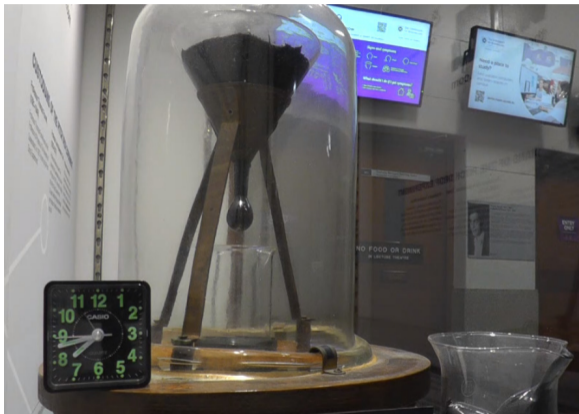
image: http://www.thetenthwatch.com/feed/

Sanity Check: Is Machine Learning the right way to solve it?

▶ Can you think of an algorithmic way to solve the problem? **No.**

Image: RobeSafe Driver Monitoring Video Dataset (RS-DMV), http://www.robesafe.com/personal/jnuevo/Datasets.html

Sanity Check: Is Machine Learning the right way to solve it?

- ▶ Can you think of an algorithmic way to solve the problem? **No.**
- ▶ It is possible to get data for the problem? **Yes.**

$\rightarrow$ machine learning sounds worth trying. ✓

Image: RobeSafe Driver Monitoring Video Dataset (RS-DMV), http://www.robesafe.com/personal/jnuevo/Datasets.html

## Example Task: Detecting Driver Fatigue

### Step 1) Decide what exactly you want

▶ **input** $x$: images (driver of a car)

▶ **output** $y \in [0, 1]$: e.g. *"how tired does the driver look?"*
        $y = 0$: totally awake, $y = 1$: sound asleep

▶ **quality measure**: e.g. $\ell(y, f(x)) = (y - f(x))^2$

▶ **model** $f_\theta$: e.g. ConvNet with certain topology, e.g. *ResNet50*

- input layer: fixed size image (scaled version of input $x$)

- output layer: single output, value $f_\theta(x) \in [0, 1]$

- parameters: $\theta$ (all weights of all layers)

▶ **goal**: find parameters $\theta^*$ such that model makes good predictions

## Step 1) Decide what exactly you want

- **input** $x$: E.g. images of the driver of a car
- **output** $y \in [0, 1]$: E.g. *"how tired does the driver look?"*
        $y = 0$: totally awake, $y = 1$: sound asleep
- **quality measure**: $\ell(\, y, f(x)\,) = (y - f(x))^2$
- **model** $f_\theta$: convolution network with certain topology
- **goal**: find parameters $\theta^*$ such that $f_{\theta^*}$ makes good predictions

### Take care that

- the outputs make sense for the inputs
    - e.g., what to output for images that don't show a person at all?

## Step 1) Decide what exactly you want

- **input** $x$: E.g. images of the driver of a car
- **output** $y \in [0, 1]$: E.g. *"how tired does the driver look?"*
  $y = 0$: totally awake, $y = 1$: sound asleep
- **quality measure**: $\ell(\, y, f(x)\,) = (y - f(x))^2$
- **model** $f_\theta$: convolution network with certain topology
- **goal**: find parameters $\theta^*$ such that $f_{\theta^*}$ makes good predictions

**Take care that**

- ▶ the outputs make sense for the inputs
  - e.g., what to output for images that don't show a person at all?
- ▶ the inputs are informative about the output you're after
  - e.g., frontal pictures? or profile? or back of the head?

## Step 1) Decide what exactly you want

- ▶ **input** $x$: E.g. images of the driver of a car
- ▶ **output** $y \in [0,1]$: E.g. *"how tired does the driver look?"*
        $y = 0$: totally awake, $y = 1$: sound asleep
- ▶ **quality measure**: $\ell(\, y, f(x)\,) = (y - f(x))^2$
- ▶ **model** $f_\theta$: convolution network with certain topology
- ▶ **goal**: find parameters $\theta^*$ such that $f_{\theta^*}$ makes good predictions

**Take care that**

- ▶ the outputs make sense for the inputs
    - – e.g., what to output for images that don't show a person at all?
- ▶ the inputs are informative about the output you're after
    - – e.g., frontal pictures? or profile? or back of the head?
- ▶ the quality measure makes sense for the task
    - – 'fatigue' (real-valued): regression, e.g. $\ell(y, f(x)) = (y - f(x))^2$
    - – 'driver identification': classification, e.g. $\ell(y, f(x)) = [\![y \neq f(x)]\!]$
    - – 'failure probability', e.g. $\ell(y, f(x)) = y \log f(x) + (1 - y) \log(1 - f(x))$

## Step 1) Decide what <u>exactly</u> you want

- ▶ **input** $x$: E.g. images of the driver of a car
- ▶ **output** $y \in [0, 1]$: E.g. *"how tired does the driver look?"*
    $y = 0$: totally awake, $y = 1$: sound asleep
- ▶ **quality measure**: $\ell(\, y, f(x)\,) = (y - f(x))^2$
- ▶ **model** $f_\theta$: convolution network with certain topology
- ▶ **goal**: find parameters $\theta^*$ such that $f_{\theta^*}$ makes good predictions

### Take care that

- ▶ the outputs make sense for the inputs
    - e.g., what to output for images that don't show a person at all?
- ▶ the inputs are informative about the output you're after
    - e.g., frontal pictures? or profile? or back of the head?
- ▶ the quality measure makes sense for the task
    - 'fatigue' (real-valued): regression, e.g. $\ell(y, f(x)) = (y - f(x))^2$
    - 'driver identification': classification, e.g. $\ell(y, f(x)) = [\![y \neq f(x)]\!]$
    - 'failure probability', e.g. $\ell(y, f(x)) = y \log f(x) + (1 - y) \log(1 - f(x))$
- ▶ the model class makes sense (later...)

## Step 2) Collect and annotate data

▶ **collect data**: images $x_1, x_2, \ldots$
▶ **annotate data**: have an expert assign 'correct' outputs: $y_1, y_2, \ldots$

**In data collection, take care that**

▶ the data reflects the situation of interest well
  – same conditions (resolution, perspective, lighting) as from car camera, . . .

▶ the data covers all situations of interest
  – drivers of all genders and ethnicities
  – diverse set of car models

▶ you have enough data
  – the more the better, but quantity is not a replacement for quality

For real-world projects, often the majority of time is invested in this step.

## Step 2) Collect and annotate data

- ▶ **collect data**: images $x_1, x_2, \ldots$
- ▶ **annotate data**: have an expert assign 'correct' outputs: $y_1, y_2, \ldots$

**For data annotation, take care that**

- ▶ the annotation is of high quality (i.e. exactly what you want)
    - – perceived 'fatigue' might be subjective
    - – how tired is '0.7' compared to '0.6' ?
    - – if multiple annotators are involved, common standards are required
    - – beware of lazy or tired annotators,
    - – beware of conventions, conversion or data entering errors

- ▶ you have enough annotated data
    - – the more the better, but quantity is not a replacement for quality

For real-world projects, this step can incur the majority of cost.

## Step 3) Model training

▶ **take a training set**, $S = \{(x_1, y_1), \ldots, (x_n, y_n)\}$,

▶ **solve the following optimization problem** to find $\theta^*$

$$\min_\theta J(\theta) \quad \text{with} \quad J(\theta) = \sum_{i=1}^{n} \ell(y_i, f_\theta(x_i))$$

## Step 3) Model training

- **take a training set**, $S = \{(x_1, y_1), \ldots, (x_n, y_n)\}$,
- **solve the following optimization problem** to find $\theta^*$

$$\min_\theta J(\theta) \quad \text{with} \quad J(\theta) = \sum_{i=1}^{n} \ell(y_i, f_\theta(x_i))$$

**What to take care of?** That's a much longer story... we'll need:

- ▶ Embrace probabilities
- ▶ Empirical risk minimization
- ▶ Overfitting / underfitting
- ▶ Regularization
- ▶ Model selection

# Embrace probabilities

## Embrace probabilities

**Most quantities in computer vision are not fully deterministic.**

- ▶ true randomness of events
  - a photon reaches a camera's CCD chip, is it detected or not?
    it depends on quantum effects, which are stochastic.

- ▶ measurement error
  - depth sensor may only be accurate $\pm 5\%$

- ▶ incomplete knowledge
  - what will be the next picture I take with my smartphone?

- ▶ insufficient representation
  - from what material is that green object made?
    hardly possible to tell from RGB, maybe easier in hyperspectral image,
    but ultimately image information is not enough

In practice, there is no difference between these!

**Probability theory allows us to deal with this.**

# Empirical Risk Minimization

## What do we want?

We have chosen (or we are given):

- ▶ general setup: inputs $\mathcal{X}$, outputs $\mathcal{Y}$
- ▶ set of models: $f_\theta$ for $\theta \in \Theta$
- ▶ loss function to judge model quality: $\ell$
- ▶ a training set $S = \{(x_1, y_1), \ldots, (x_n, y_n)\}$

## What do we want?

## What do we want?

We have chosen (or we are given):

- ▶ general setup: inputs $\mathcal{X}$, outputs $\mathcal{Y}$
- ▶ set of models: $f_\theta$ for $\theta \in \Theta$
- ▶ loss function to judge model quality: $\ell$
- ▶ a training set $S = \{(x_1, y_1), \ldots, (x_n, y_n)\}$

### What do we want?

A model $f_{\theta^*}$ that works well (=has small loss) when applied to the data in $S$.

## What do we want?

We have chosen (or we are given):

- ▶ general setup: inputs $\mathcal{X}$, outputs $\mathcal{Y}$
- ▶ set of models: $f_\theta$ for $\theta \in \Theta$
- ▶ loss function to judge model quality: $\ell$
- ▶ a training set $S = \{(x_1, y_1), \ldots, (x_n, y_n)\}$

### What do we want?

A model $f_{\theta^*}$ that works well (=has small loss) when applied to the data in $S$.

**Why?** If you are interested in the label of some $x_i$, just look it up in $S$.

## What do we want?

We have chosen (or we are given):

- ▶ general setup: inputs $\mathcal{X}$, outputs $\mathcal{Y}$
- ▶ set of models: $f_\theta$ for $\theta \in \Theta$
- ▶ loss function to judge model quality: $\ell$
- ▶ a training set $S = \{(x_1, y_1), \ldots, (x_n, y_n)\}$

## What do we want?

A model $f_{\theta^*}$ that works well (=has small loss) when applied to the data in $S$.

**Why?** If you are interested in the label of some $x_i$, just look it up in $S$.

## What do we (really) want?

## What do we want?

We have chosen (or we are given):

- ▶ general setup: inputs $\mathcal{X}$, outputs $\mathcal{Y}$
- ▶ set of models: $f_\theta$ for $\theta \in \Theta$
- ▶ loss function to judge model quality: $\ell$
- ▶ a training set $S = \{(x_1, y_1), \ldots, (x_n, y_n)\}$

### What do we want?

A model $f_{\theta^*}$ that works well (=has small loss) when applied to the data in $S$.

**Why?** If you are interested in the label of some $x_i$, just look it up in $S$.

### What do we (really) want?

Let $S' = \{(x'_1, y'_1), \ldots, (x'_n, y'_n)\}$ be a test set. We want a model $\theta^*$ that $f_{\theta^*}$ that works well (=has small loss) when applied to the data in $S'$.

## What do we want?

We have chosen (or we are given):

- ▶ general setup: inputs $\mathcal{X}$, outputs $\mathcal{Y}$
- ▶ set of models: $f_\theta$ for $\theta \in \Theta$
- ▶ loss function to judge model quality: $\ell$
- ▶ a training set $S = \{(x_1, y_1), \ldots, (x_n, y_n)\}$

### What do we want?

A model $f_{\theta^*}$ that works well (=has small loss) when applied to the data in $S$.

**Why?** If you are interested in the label of some $x_i$, just look it up in $S$.

### What do we (really) want?

Let $S' = \{(x'_1, y'_1), \ldots, (x'_n, y'_n)\}$ be a test set. We want a model $\theta^*$ that $f_{\theta^*}$ that works well (=has small loss) when applied to the data in $S'$.

**Still, why?** If you are interested in the label of some $x'_i$, just look it up in $S'$.

## What do we (really, really) want?

A model $f_\theta$ that works well (=has small loss) when we apply it to future data.

## What do we (really, really) want?

A model $f_\theta$ that works well (=has small loss) when we apply it to future data.

**Problem:** we don't know what inputs the future will bring!

**Probabilities to the rescue:**

- ▶ we are **uncertain** about future input data $\rightarrow$ use **random variable** $X$
- ▶ $\mathcal{X}$: all possible images, $p(x)$ probability to see any $x \in \mathcal{X}$

## What do we (really, really) want?

A model $f_\theta$ that works well (=has small loss) when we apply it to future data.

**Problem:** we don't know what inputs the future will bring!

**Probabilities to the rescue:**
- ▶ we are **uncertain** about future input data → use **random variable** $X$
- ▶ $\mathcal{X}$: all possible images, $p(x)$ probability to see any $x \in \mathcal{X}$

**Problem:** we don't know what the right outputs are for the inputs.

**Probabilities to the rescue:**
- ▶ we are **uncertain** about the outputs → use **random variable** $Y$
- ▶ $\mathcal{Y}$: all possible outputs, $p(y|x)$ probability that $y \in \mathcal{Y}$ is correct for some $x \in \mathcal{X}$

Note: we don't pretend that we know $p(x)$ or $p(y|x)$, we just assume they exist.

## What do we want formally?

A model $f_\theta$ with small expected loss (aka "risk" or "test error")

$$\mathcal{R} = \mathbb{E}_{x \sim p(x)} \mathbb{E}_{y \sim p(y|x)} [\ell(y, f_\theta(x))] = \mathbb{E}_{(x,y) \sim p(x,y)} [\ell(y, f_\theta(x))]$$

**New problem:** we can't compute $\mathcal{R}$, because we don't know $p(x, y)$ (nor $p(x)$, nor $p(y|x)$)

## What do we want formally?

A model $f_\theta$ with small expected loss (aka "risk" or "test error")

$$\mathcal{R} = \mathbb{E}_{x \sim p(x)} \mathbb{E}_{y \sim p(y|x)} [\, \ell(\, y, f_\theta(x)\,)\,] = \mathbb{E}_{(x,y) \sim p(x,y)} [\, \ell(\, y, f_\theta(x)\,)\,]$$

**New problem:** we can't compute $\mathcal{R}$, because we don't know $p(x, y)$ (nor $p(x)$, nor $p(y|x)$)

**Good news:** We can estimate it!

## Empirical risk

Let $S = \{(x_1, y_1), \ldots, (x_n, y_n)\}$ be a set of input-output pairs. Then

$$\hat{\mathcal{R}} = \frac{1}{n} \sum_{i=1}^{n} \ell(\, y_i, f_\theta(x_i))$$

is called the empirical risk of $\mathcal{R}$ w.r.t. $S$.      (short: training error)

We can easily compute $\hat{\mathcal{R}}$. Can we use it as a drop-in replacement for $\mathcal{R}$?

**Is $\hat{\mathcal{R}} = \dfrac{1}{n} \sum\limits_{i=1}^{n} \ell(\, y_i, f_\theta(x_i))$ a good estimator of $\mathcal{R}$?**

**Is $\hat{\mathcal{R}} = \dfrac{1}{n} \sum\limits_{i=1}^{n} \ell(\, y_i, f_\theta(x_i))$ a good estimator of $\mathcal{R}$? That depends on data we use.**

Independent and identically distributed (i.i.d.) training data

If the samples $x_1, \ldots, x_n$ are sampled independently from $p(x)$ and the outputs $y_1, \ldots, y_n$ are sampled independently from $p(y|x)$, then $\hat{\mathcal{R}}$ is an unbiased and consistent estimator of $\mathcal{R}$:

$$\mathbb{E}_{(x_1, y_1), \ldots, (x_n, y_n)}[\hat{\mathcal{R}}] = \mathcal{R} \qquad \text{and} \qquad \mathrm{Var}(\hat{\mathcal{R}}) \to 0 \text{ with speed } O(\tfrac{1}{n})$$

(essentially, $\hat{\mathcal{R}}$ is a noisy version of $\mathcal{R}$, and the more data we have, the smaller the noise)

**Is $\hat{\mathcal{R}} = \dfrac{1}{n} \sum_{i=1}^{n} \ell(y_i, f_\theta(x_i))$ a good estimator of $\mathcal{R}$? That depends on data we use.**

Independent and identically distributed (i.i.d.) training data

If the samples $x_1, \ldots, x_n$ are sampled independently from $p(x)$ and the outputs $y_1, \ldots, y_n$ are sampled independently from $p(y|x)$, then $\hat{\mathcal{R}}$ is an unbiased and consistent estimator of $\mathcal{R}$:

$$\mathbb{E}_{(x_1,y_1),\ldots,(x_n,y_n)}[\hat{\mathcal{R}}] = \mathcal{R} \qquad \text{and} \qquad \mathrm{Var}(\hat{\mathcal{R}}) \to 0 \;\; \text{with speed } O(\tfrac{1}{n})$$

(essentially, $\hat{\mathcal{R}}$ is a noisy version of $\mathcal{R}$, and the more data we have, the smaller the noise)

What if the data $(x_1, y_1), \ldots, (x_n, y_n)$ are not independent (e.g. video frames)?

▶ $\hat{\mathcal{R}}$ is unbiased, but $\mathrm{Var}[\hat{\mathcal{R}}]$ will converge slower to $0$ (potentially very very slowly)

What if the distribution of the samples $(x_1, y_1), \ldots, (x_n, y_n)$ is different from $p$ (biased inputs, biased outputs or both)?

▶ $\hat{\mathcal{R}}$ can differ a lot from $\mathcal{R}$ ($\to$ "domain adaptation", "out-of-distribution", "fairness")

## Step 2) Collect and annotate data

- ▶ collect examples: $x_1, x_2, \ldots$
- ▶ have an expert annotate them with 'correct' outputs: $y_1, y_2, \ldots$

**Take care that:**

- ▶ data comes from the data distribution we actually care about ("prediction time")
- ▶ examples are independent
- ▶ output annotation is a good as possible

## Step 2) Collect and annotate data

- ▶ collect examples: $x_1, x_2, \ldots$
- ▶ have an expert annotate them with 'correct' outputs: $y_1, y_2, \ldots$

**Take care that:**

- ▶ data comes from the data distribution we actually care about ("prediction time")
- ▶ examples are independent
- ▶ output annotation is a good as possible

## Step 3) Model training

Take a training set, $S = \{(x_1, y_1), \ldots, (x_n, y_n)\}$, and find $\theta^*$ by minimizing the loss on the training set:

$$\min_{\theta} J(\theta) \quad \text{with} \quad J(\theta) = \sum_{i=1}^{n} \ell(y_i, f_\theta(x_i))$$

Observation: $J(\theta)$ is simply the *empirical risk* $\hat{\mathcal{R}}$ (up to an irrelevant factor $\frac{1}{n}$)

Learning by **Empirical Risk Minimization**

## Summary

Learning is about finding a model that works on future data.

The true quantity of interest is the expected error on future data:

$$\mathcal{R} = \mathbb{E}_{(x,y) \sim p(x,y)} \ell(y, f(x)) \qquad \text{(test error)}$$

We cannot compute $\mathcal{R}$, but we can estimate it.

For a training set, $S = \{(x_1, y_1), \ldots, (x_n, y_n)\}$, we can compute the average error:

$$\hat{\mathcal{R}} = \sum_{i=1}^{n} \ell(y_i, f(x_i)) \qquad \text{(training error)}$$

Training data should ideally be i.i.d. from the right distribution

If the training data are sampled independently and all from the distribution $p(x, y)$ then $\hat{\mathcal{R}}$ is a good (unbiased and consistent) estimator of $\mathcal{R}$.

Learning $\equiv$ empirical risk minimization

The core of most learning methods is to minimize the training error.

# Overfitting / underfitting

We found a model $f_{\theta^*}$ by minimizing the training error $\hat{\mathcal{R}}$.

Q: Will it work well on future data, i.e. have small test error, $\mathcal{R}$?

A: **Unfortunately, that is not guaranteed.**

We found a model $f_{\theta^*}$ by minimizing the training error $\hat{\mathcal{R}}$.

Q: Will it work well on future data, i.e. have small test error, $\mathcal{R}$?

A: **Unfortunately, that is not guaranteed.**

## Relation between training error and test error

Example: 1D curve fitting



training points

We found a model $f_{\theta^*}$ by minimizing the training error $\hat{\mathcal{R}}$.

Q: Will it work well on future data, i.e. have small test error, $\mathcal{R}$?

A: **Unfortunately, that is not guaranteed.**

## Relation between training error and test error

Example: 1D curve fitting



- degree 2 fit, $\hat{\mathcal{R}} = 8.44$
- true signal $\mathcal{R} = 14.64$
- training points

best learned polynomial of degree 2: large $\hat{\mathcal{R}}$, large $\mathcal{R}$

We found a model $f_{\theta^*}$ by minimizing the training error $\hat{\mathcal{R}}$.

Q: Will it work well on future data, i.e. have small test error, $\mathcal{R}$?

A: **Unfortunately, that is not guaranteed.**

Relation between training error and test error

Example: 1D curve fitting



best learned polynomial of degree 7: small $\hat{\mathcal{R}}$, small $\mathcal{R}$

We found a model $f_{\theta^*}$ by minimizing the training error $\hat{\mathcal{R}}$.

Q: Will it work well on future data, i.e. have small test error, $\mathcal{R}$?

A: **Unfortunately, that is not guaranteed.**

### Relation between training error and test error

Example: 1D curve fitting



best learned polynomial of degree 12: small $\hat{\mathcal{R}}$, large $\mathcal{R}$

We found a model $f_{\theta^*}$ by minimizing the training error $\hat{\mathcal{R}}$.

Q: Will its test error, $\mathcal{R}$, be small?

A: **Unfortunately, that is not guaranteed.**

## Underfitting/Overfitting



| | |
|---|---|
| **Underfitting** | **Overfitting** |
| (to some extent) detectable from $\hat{\mathcal{R}}$ | not detectable from $\hat{\mathcal{R}}$ ! |

Choosing a model based on $\hat{\mathcal{R}}$ vs. $\mathcal{R}$

test error $\mathcal{R}$ for 7 different models/parameter choices

test error $\mathcal{R}$ for 7 different models/parameter choices

Choosing a model based on $\hat{\mathcal{R}}$ vs. $\mathcal{R}$

Legend:
- $\mathcal{R}(\theta_i)$
- $\hat{\mathcal{R}}_{S_3}(\theta_i)$

x-axis: $\theta_i$

training error $\hat{\mathcal{R}}$ for a training set, $S$

training errors $\hat{\mathcal{R}}$ for $5$ possible training sets

Choosing hypothesis based on $\hat{\mathcal{R}}$ vs. $\mathcal{R}$

Legend:
- $\mathcal{R}(\theta_i)$
- $\hat{\mathcal{R}}_{S_1}(\theta_i)$
- $\hat{\mathcal{R}}_{S_2}(\theta_i)$
- $\hat{\mathcal{R}}_{S_3}(\theta_i)$
- $\hat{\mathcal{R}}_{S_4}(\theta_i)$
- $\hat{\mathcal{R}}_{S_5}(\theta_i)$

model with smallest training error can have high test error

**The training error does not tell us how good a model really is. So, what does?**

**The training error does not tell us how good a model really is. So, what does?**

## Step 4) Model evaluation

Take new data, $S' = \{(x'_1, y'_1), \ldots, (x'_m, y'_m)\}$, and compute the model performance as

$$\hat{\mathcal{R}}_{\text{tst}} = \frac{1}{m} \sum_{i=1}^{m} \ell(\, y'_i, f_{\theta^*}(x'_i)\,)$$

**Take care that:**

▶ data and annotation for evaluation comes from the real distribution

▶ examples are independent from each other

▶ data is independent of the chosen model,
   in particular: model architecture and parameter depend in no way on $S'$

If all of these are fulfilled:

▶ $\hat{\mathcal{R}}_{\text{tst}}$ is an unbiased estimate of $\mathcal{R}$, and we can expect $\hat{\mathcal{R}}_{\text{tst}} = \mathcal{R} + O(\frac{1}{\sqrt{m}})$

**The training error does not tell us how good a model really is. So, what does?**

## Step 4) Model evaluation

Take new data, $S' = \{(x'_1, y'_1), \ldots, (x'_m, y'_m)\}$, and compute the model performance as

$$\hat{\mathcal{R}}_{\text{tst}} = \frac{1}{m} \sum_{i=1}^{m} \ell(\, y'_i, f_{\theta^*}(x'_i)\,)$$

**Take care that:**

▶ data and annotation for evaluation comes from the real distribution
▶ examples are independent from each other
▶ data is independent of the chosen model,
  in particular: model architecture and parameter depend in no way on $S'$

If all of these are fulfilled:

▶ $\hat{\mathcal{R}}_{\text{tst}}$ is an unbiased estimate of $\mathcal{R}$, and we can expect $\hat{\mathcal{R}}_{\text{tst}} = \mathcal{R} + O(\frac{1}{\sqrt{m}})$

**The training error does not tell us how good a model really is. So, what does?**

## Step 4) Model evaluation

Take new data, $S' = \{(x'_1, y'_1), \ldots, (x'_m, y'_m)\}$, and compute the model performance as

$$\hat{\mathcal{R}}_{\text{tst}} = \frac{1}{m} \sum_{i=1}^{m} \ell(\, y'_i, f_{\theta^*}(x'_i)\,)$$

**Take care that:**

▶ data and annotation for evaluation comes from the real distribution

▶ examples are independent from each other

▶ data is independent of the chosen model,
  in particular: model architecture and parameter depend in no way on $S'$

If all of these are fulfilled:

▶ $\hat{\mathcal{R}}_{\text{tst}}$ is an unbiased estimate of $\mathcal{R}$, and we can expect $\hat{\mathcal{R}}_{\text{tst}} = \mathcal{R} + O(\frac{1}{\sqrt{m}})$

Task: **Driver fatigue**. Available training data are images from many driving videos.

Which of these are proper **test sets**?

Task: **Driver fatigue**. Available training data are images from many driving videos.

Which of these are proper **test sets**?

▶ download more videos of the same drivers

Task: **Driver fatigue**. Available training data are images from many driving videos.

Which of these are proper **test sets**?

▶ download more videos of the same drivers ✗

Task: **Driver fatigue**. Available training data are images from many driving videos.

Which of these are proper **test sets**?

▶ download more videos of the same drivers ✗

▶ download more videos of the same drivers, but make sure that drivers who were asleep in the training set only appear as awake in the test set and vice versa

Task: **Driver fatigue**. Available training data are images from many driving videos.

Which of these are proper **test sets**?

▶ download more videos of the same drivers ✗

▶ download more videos of the same drivers, but make sure that drivers who were asleep in the training set only appear as awake in the test set and vice versa ✗

Task: **Driver fatigue**. Available training data are images from many driving videos.

Which of these are proper **test sets**?

► download more videos of the same drivers ✗

► download more videos of the same drivers, but make sure that drivers who were asleep in the training set only appear as awake in the test set and vice versa ✗

► download more videos of different drivers

Task: **Driver fatigue**. Available training data are images from many driving videos.

Which of these are proper **test sets**?

▶ download more videos of the same drivers ✗

▶ download more videos of the same drivers, but make sure that drivers who were asleep in the training set only appear as awake in the test set and vice versa ✗

▶ download more videos of different drivers ✓

Task: **Driver fatigue**. Available training data are images from many driving videos.

Which of these are proper **test sets**?

▶ download more videos of the same drivers ✗

▶ download more videos of the same drivers, but make sure that drivers who were asleep in the training set only appear as awake in the test set and vice versa ✗

▶ download more videos of different drivers ✓

▶ download videos recorded in winter whereas the training set was recorded in summer

Task: **Driver fatigue**. Available training data are images from many driving videos.

Which of these are proper **test sets**?

▶ download more videos of the same drivers ✗

▶ download more videos of the same drivers, but make sure that drivers who were asleep in the training set only appear as awake in the test set and vice versa ✗

▶ download more videos of different drivers ✓

▶ download videos recorded in winter whereas the training set was recorded in summer ✗

Task: **Driver fatigue**. Available training data are images from many driving videos.

Which of these are proper **test sets**?

▶ download more videos of the same drivers ✗

▶ download more videos of the same drivers, but make sure that drivers who were asleep in the training set only appear as awake in the test set and vice versa ✗

▶ download more videos of different drivers ✓

▶ download videos recorded in winter whereas the training set was recorded in summer ✗

▶ download videos from a different video platform than the training set

Task: **Driver fatigue**. Available training data are images from many driving videos.

Which of these are proper **test sets**?

▶ download more videos of the same drivers ✗

▶ download more videos of the same drivers, but make sure that drivers who were asleep in the training set only appear as awake in the test set and vice versa ✗

▶ download more videos of different drivers ✓

▶ download videos recorded in winter whereas the training set was recorded in summer ✗

▶ download videos from a different video platform than the training set ✗

Task: **Driver fatigue**. Available training data are images from many driving videos.

Which of these are proper **test sets**?

- ▶ download more videos of the same drivers ✗

- ▶ download more videos of the same drivers, but make sure that drivers who were asleep in the training set only appear as awake in the test set and vice versa ✗

- ▶ download more videos of different drivers ✓

- ▶ download videos recorded in winter whereas the training set was recorded in summer ✗

- ▶ download videos from a different video platform than the training set ✗

**Summary: It's not always obvious how to obtain a good test set.**

You are given a **test set**. What are you allowed to use it for?

You are given a **test set**. What are you allowed to use it for?

- ▶ train the model

You are given a **test set**. What are you allowed to use it for?

- train the model ✗

You are given a **test set**. What are you allowed to use it for?

- ▶ train the model ✗
- ▶ decide when to stop training (e.g. when test error does not change anymore)

You are given a **test set**. What are you allowed to use it for?

▶ train the model ✗

▶ decide when to stop training (e.g. when test error does not change anymore) ✗

You are given a **test set**. What are you allowed to use it for?

▶ train the model ✗

▶ decide when to stop training (e.g. when test error does not change anymore) ✗

▶ adjust hyper-parameters (e.g. learning rate)

You are given a **test set**. What are you allowed to use it for?

- ▶ train the model ✗
- ▶ decide when to stop training (e.g. when test error does not change anymore) ✗
- ▶ adjust hyper-parameters (e.g. learning rate) ✗

You are given a **test set**. What are you allowed to use it for?

▶ train the model ✗

▶ decide when to stop training (e.g. when test error does not change anymore) ✗

▶ adjust hyper-parameters (e.g. learning rate) ✗

▶ get upset, trash your model, and switch to a completely different network architecture

You are given a **test set**. What are you allowed to use it for?

- ▶ train the model ✗
- ▶ decide when to stop training (e.g. when test error does not change anymore) ✗
- ▶ adjust hyper-parameters (e.g. learning rate) ✗
- ▶ get upset, trash your model, and switch to a completely different network architecture ✗

You are given a **test set**. What are you allowed to use it for?

- ▶ train the model ✗
- ▶ decide when to stop training (e.g. when test error does not change anymore) ✗
- ▶ adjust hyper-parameters (e.g. learning rate) ✗
- ▶ get upset, trash your model, and switch to a completely different network architecture ✗

**"overfitting to the test set"**

- ▶ evaluate your final model's accuracy in order to compare it to the literature

You are given a **test set**. What are you allowed to use it for?

- ▶ train the model ✗
- ▶ decide when to stop training (e.g. when test error does not change anymore) ✗
- ▶ adjust hyper-parameters (e.g. learning rate) ✗
- ▶ get upset, trash your model, and switch to a completely different network architecture ✗

**"overfitting to the test set"**

- ▶ evaluate your final model's accuracy in order to compare it to the literature ✓

**Avoid overfitting to the test set!** It

- ▶ invalidates your results (you'll think your model is better than it is),
- ▶ undermines the credibility of your experimental evaluation.

Test sets are precious! They can only be used once!

Trap 1: additional hyper-parameters

Imagine you propose a new term, $\Gamma(\theta)$ for a loss function

▶ previous work:

$$\min_\theta \quad F(\theta)$$

▶ proposed method:

$$\min_\theta \quad F(\theta) + \lambda\Gamma(\theta)$$

▶ hyperparameters: $\lambda \in \{0, 0.1, \ldots, 0.9, 1\}$

Trap 1: additional hyper-parameters

Imagine you propose a new term, $\Gamma(\theta)$ for a loss function

▶ previous work:

$$\min_{\theta} \quad F(\theta)$$

▶ proposed method:

$$\min_{\theta} \quad F(\theta) + \lambda\Gamma(\theta)$$

▶ hyperparameters: $\lambda \in \{0, 0.1, \ldots, 0.9, 1\}$

If you select $\lambda$ using the test set, results for the new method will always be at least as good as for the old one, regardless what $\Gamma$ was. $\rightarrow$ experimental results prove nothing!

Overfitting to the test set can be avoided by not making the test data public.

Example: ChaLearn Connectomics Challenge

▶ We are given labeled training data.

▶ We build/train a model.

▶ We upload the model in executable form to a server.

▶ The server applies the model to test data and evaluates the results.

Example: ImageNet ILSVR Challenge

▶ We are given labeled training data and unlabeled test data.

▶ We build/train a model.

▶ We apply the model to the test data.

▶ We upload the model predictions to a server.

▶ The server evaluates the results.

# Regularization

**overfitting**

$\hat{\mathcal{R}}$ **vs.** $\mathcal{R}$

**How can we prevent overfitting when learning a model?**

Choosing hypothesis based on $\hat{\mathcal{R}}$ vs. $\mathcal{R}$

Legend:
- $\mathcal{R}(\theta_i)$
- $\hat{\mathcal{R}}_{S_1}(\theta_i)$
- $\hat{\mathcal{R}}_{S_2}(\theta_i)$
- $\hat{\mathcal{R}}_{S_3}(\theta_i)$
- $\hat{\mathcal{R}}_{S_4}(\theta_i)$
- $\hat{\mathcal{R}}_{S_5}(\theta_i)$

Choosing hypothesis based on $\hat{\mathcal{R}}$ vs. $\mathcal{R}$

Legend:
- $\mathcal{R}(\theta_i)$
- $\hat{\mathcal{R}}_{S_1}(\theta_i)$
- $\hat{\mathcal{R}}_{S_2}(\theta_i)$
- $\hat{\mathcal{R}}_{S_3}(\theta_i)$
- $\hat{\mathcal{R}}_{S_4}(\theta_i)$
- $\hat{\mathcal{R}}_{S_5}(\theta_i)$

larger training set $\rightarrow$ smaller variance of $\hat{\mathcal{R}}$

Choosing hypothesis based on $\hat{\mathcal{R}}$ vs. $\mathcal{R}$

lower probability that $\hat{\mathcal{R}}$ differs strongly from $\mathcal{R}$ → less overfitting

Choosing hypothesis based on $\hat{\mathcal{R}}$ vs. $\mathcal{R}$

Legend:
- $\mathcal{R}(\theta_i)$
- $\hat{\mathcal{R}}_{S_1}(\theta_i)$
- $\hat{\mathcal{R}}_{S_2}(\theta_i)$
- $\hat{\mathcal{R}}_{S_3}(\theta_i)$
- $\hat{\mathcal{R}}_{S_4}(\theta_i)$
- $\hat{\mathcal{R}}_{S_5}(\theta_i)$

$\theta_i$

Choosing hypothesis based on $\hat{\mathcal{R}}$ vs. $\mathcal{R}$

Legend:
- $\mathcal{R}(\theta_i)$
- $\hat{\mathcal{R}}_{S_1}(\theta_i)$
- $\hat{\mathcal{R}}_{S_2}(\theta_i)$
- $\hat{\mathcal{R}}_{S_3}(\theta_i)$
- $\hat{\mathcal{R}}_{S_4}(\theta_i)$
- $\hat{\mathcal{R}}_{S_5}(\theta_i)$

fewer models $\rightarrow$ lower probability of a model with small $\hat{\mathcal{R}}$s but high $\mathcal{R}$

Choosing hypothesis based on $\hat{\mathcal{R}}$ vs. $\mathcal{R}$

Legend:
- $\mathcal{R}(\theta_i)$
- $\hat{\mathcal{R}}_{S_1}(\theta_i)$
- $\hat{\mathcal{R}}_{S_2}(\theta_i)$
- $\hat{\mathcal{R}}_{S_3}(\theta_i)$
- $\hat{\mathcal{R}}_{S_4}(\theta_i)$
- $\hat{\mathcal{R}}_{S_5}(\theta_i)$

fewer models $\rightarrow$ lower probability of a model with small $\hat{\mathcal{R}}$s but high $\mathcal{R}$

to few models select to from $\rightarrow$ danger that no model with low $\mathcal{R}$ is left!

Choosing hypothesis based on $\hat{\mathcal{R}}$ vs. $\mathcal{R}$

$\mathcal{R}(\theta_i)$

$\theta_i$

to few models select to from $\rightarrow$ danger that no model with low $\mathcal{R}$ is left!

Choosing hypothesis based on $\hat{\mathcal{R}}$ vs. $\mathcal{R}$

Legend:
- $\mathcal{R}(\theta_i)$
- $\hat{\mathcal{R}}_{S_1}(\theta_i)$
- $\hat{\mathcal{R}}_{S_2}(\theta_i)$
- $\hat{\mathcal{R}}_{S_3}(\theta_i)$
- $\hat{\mathcal{R}}_{S_4}(\theta_i)$
- $\hat{\mathcal{R}}_{S_5}(\theta_i)$

**Underfitting!**

Choosing hypothesis based on $\hat{\mathcal{R}}$ vs. $\mathcal{R}$

**Underfitting!**

Models with big difference between training error and test error are typically **extreme cases**:

- ▶ a large number of model parameters
- ▶ model parameters take large values
- ▶ learned function is very non-smooth



coeffs: $\theta_i \in [-2.4, 4.6]$

coeffs: $\theta_i \in [-1312.5, 1136.6]$

Models with big difference between training error and test error are typically **extreme cases**:

- ▶ a large number of model parameters
- ▶ model parameters take large values
- ▶ learned function is very non-smooth



coeffs: $\theta_i \in [-2.4, 4.6]$        coeffs: $\theta_i \in [-1312.5, 1136.6]$

**Regularization:** avoid overfitting by preventing extremes to occur

- ▶ explicit regularization (changing the objective function)
- ▶ implicit regularization (modifying the optimization procedure)

## Explicit regularization

Add a regularization term (=regularizer) to the empirical risk that penalizes extreme parameter values.

### Regularized risk minimization

Take a training set, $S = \{(x_1, y_1), \ldots, (x_n, y_n)\}$, find $\theta^*$ by solving,

$$\min_\theta J_\lambda(\theta) \quad \text{with} \quad J_\lambda(\theta) = \underbrace{\sum_{i=1}^n \ell(y_i, f_\theta(x_i))}_{\text{empirical risk}} + \underbrace{\lambda\Omega(\theta)}_{\text{regularizer}}$$

e.g. with $\quad \Omega(\theta) = \|\theta\|_{L^2}^2 = \sum_j \theta_j^2 \quad$ or $\quad \Omega(\theta) = \|\theta\|_{L^1} = \sum_j |\theta_j|$

## Explicit regularization

Add a regularization term (=regularizer) to the empirical risk that penalizes extreme parameter values.

### Regularized risk minimization

Take a training set, $S = \{(x_1, y_1), \ldots, (x_n, y_n)\}$, find $\theta^*$ by solving,

$$\min_\theta J_\lambda(\theta) \quad \text{with} \quad J_\lambda(\theta) = \underbrace{\sum_{i=1}^n \ell(y_i, f_\theta(x_i))}_{\text{empirical risk}} + \underbrace{\lambda \Omega(\theta)}_{\text{regularizer}}$$

$$\text{e.g. with} \quad \Omega(\theta) = \|\theta\|_{L^2}^2 = \sum_j \theta_j^2 \quad \text{or} \quad \Omega(\theta) = \|\theta\|_{L^1} = \sum_j |\theta_j|$$

Optimization searches model with small training error <u>and</u> small parameter values.

**Regularization (hyper)parameter** $\lambda \geq 0$: trade-off between both.
- ▶ $\lambda = 0$: empirical risk minimization (risk of overfitting)
- ▶ $\lambda \to \infty$: all parameters $0$ (risk of underfitting)

$$\min_w J_\lambda(w) \quad \text{for} \quad J_\lambda(w) = \sum_{i=1}^n (w^\top x_i - y_i)^2 + \lambda\|w\|^2$$

$$\min_w J_\lambda(w) \quad \text{for} \quad J_\lambda(w) = \sum_{i=1}^n (w^\top x_i - y_i)^2 + \lambda \|w\|^2$$

Train/val error for minimizer of $J_\lambda$ with varying amounts of regularization:



eye dataset: 737 examples for training, 736 examples for evaluation

$$\min_w J_\lambda(w) \quad \text{for} \quad J_\lambda(w) = \sum_{i=1}^{n} (w^\top x_i - y_i)^2 + \lambda \|w\|^2$$

Train/val error for minimizer of $J_\lambda$ with varying amounts of regularization:



eye dataset: 737 examples for training, 736 examples for evaluation

$$\min_w J_\lambda(w) \quad \text{for} \quad J_\lambda(w) = \sum_{i=1}^{n} (w^\top x_i - y_i)^2 + \lambda \|w\|^2$$

Train/val error for minimizer of $J_\lambda$ with varying amounts of regularization:



eye dataset: 737 examples for training, 736 examples for evaluation

Training error, $\hat{\mathcal{R}}$, is a noisy estimate of the test error, $\mathcal{R}$

- original risk $\hat{\mathcal{R}}$ is unbiased, but variance can be huge
- regularization introduces a bias, but reduces variance
- for $\lambda \to \infty$, the variance goes to $0$, but the bias gets very big



Image: adapted from http://scott.fortmann-roe.com/docs/BiasVariance.html

Training error, $\hat{\mathcal{R}}$, is a noisy estimate of the test error, $\mathcal{R}$

- original risk $\hat{\mathcal{R}}$ is unbiased, but variance can be huge
- regularization introduces a bias, but reduces variance
- for $\lambda \to \infty$, the variance goes to $0$, but the bias gets very big



Image: adapted from http://scott.fortmann-roe.com/docs/BiasVariance.html

Training error, $\hat{\mathcal{R}}$, is a noisy estimate of the test error, $\mathcal{R}$

- original risk $\hat{\mathcal{R}}$ is unbiased, but variance can be huge
- regularization introduces a bias, but reduces variance
- for $\lambda \to \infty$, the variance goes to $0$, but the bias gets very big



Image: adapted from http://scott.fortmann-roe.com/docs/BiasVariance.html

Training error, $\hat{\mathcal{R}}$, is a noisy estimate of the test error, $\mathcal{R}$

- original risk $\hat{\mathcal{R}}$ is unbiased, but variance can be huge
- regularization introduces a bias, but reduces variance
- for $\lambda \to \infty$, the variance goes to $0$, but the bias gets very big



Image: adapted from http://scott.fortmann-roe.com/docs/BiasVariance.html

## Implicit regularization

Numerical optimization is performed iteratively, e.g. gradient descent

### Gradient descent optimization

- initialize $\theta^{(0)}$
- **for** $t = 1, 2, \ldots$
- $\quad \theta^{(t)} \leftarrow \theta^{(t-1)} - \eta_t \nabla_\theta J(\theta^{(t-1)})$ $\qquad (\eta_t \in \mathbb{R}$ is some stepsize rule$)$
- **until convergence**

**Implicit regularization** prevent overfitting by modifying these steps, e.g.

- early stopping
- weight decay
- data augmentation
- dropout
- inductive bias $\qquad \longrightarrow$ **Intro to Deep Learning**

# Model Selection

**Which of model class?**

▶ linear classifier, gradient boosted decision tree, neural network?

▶ ResNet? DenseNet? EfficientNet? Transformer? MLPMixer? GAN? Diffusion Network?

▶ which layers? how many? how wide?

**Which data representation?**

▶ raw pixels? manual preprocessing? pretrained features?

**Which hyperparameters?**

▶ batchsize? learning rate? number of epochs? which optimizer?

▶ regularization or not? which one? how much?

**Which of model class?**

- ▶ linear classifier, gradient boosted decision tree, neural network?
- ▶ ResNet? DenseNet? EfficientNet? Transformer? MLPMixer? GAN? Diffusion Network?
- ▶ which layers? how many? how wide?

**Which data representation?**

- ▶ raw pixels? manual preprocessing? pretrained features?

**Which hyperparameters?**

- ▶ batchsize? learning rate? number of epochs? which optimizer?
- ▶ regularization or not? which one? how much?

### Model Selection

**Which of model class?**
- ▶ linear classifier, gradient boosted decision tree, neural network?
- ▶ ResNet? DenseNet? EfficientNet? Transformer? MLPMixer? GAN? Diffusion Network?
- ▶ which layers? how many? how wide?

**Which data representation?**
- ▶ raw pixels? manual preprocessing? pretrained features?

**Which hyperparameters?**
- ▶ batchsize? learning rate? number of epochs? which optimizer?
- ▶ regularization or not? which one? how much?

## Model Selection

Model selection is a difficult (and often underestimated) problem:
- ▶ we can't decide based on training error: we won't catch overfitting
- ▶ we can't decide based on test error: if we use test data to select (hyper)parameters, we're overfitting to the test set and the result is not a good estimate of true risk anymore

**Problem:** We want to evaluate a model on different data than the data it was trained on, but not the test set.

**Solution:** Emulate the train/test split from only the data available for training.

- full data ("**trainval**") $\longrightarrow$ **train** and **validation** sets

- aim for train/validation sets to have same relation as train/test sets:
  – same data distribution
  – trained model independent of **validation** data
  – both sets as big as possible

**Driver fatigue:** Available data are images from many driving videos of some group of people

Which of these are good **train/validation splits**?

**Driver fatigue:** Available data are images from many driving videos of some group of people

Which of these are good **train/validation splits**?

► split the images randomly into two subsets, one for training and the other for validation

**Driver fatigue:** Available data are images from many driving videos of some group of people

Which of these are good **train/validation splits**?

▶ split the images randomly into two subsets, one for training and the other for validation ✗

**Driver fatigue:** Available data are images from many driving videos of some group of people

Which of these are good **train/validation splits**?

▶ split the images randomly into two subsets, one for training and the other for validation ✗

▶ split the videos randomly into two subsets. . .

**Driver fatigue:** Available data are images from many driving videos of some group of people

Which of these are good **train/validation splits**?

- ▶ split the images randomly into two subsets, one for training and the other for validation ✗
- ▶ split the videos randomly into two subsets... ?  → what about repeated drivers?

**Driver fatigue:** Available data are images from many driving videos of some group of people

Which of these are good **train/validation splits**?

- ▶ split the images randomly into two subsets, one for training and the other for validation ✗
- ▶ split the videos randomly into two subsets. . .   **?**   → what about repeated drivers?
- ▶ split the drivers randomly into two subsets. . .

**Driver fatigue:** Available data are images from many driving videos of some group of people

Which of these are good **train/validation splits**?

- ▶ split the images randomly into two subsets, one for training and the other for validation ✗
- ▶ split the videos randomly into two subsets... ? → what about repeated drivers?
- ▶ split the drivers randomly into two subsets... ? → repeated car models?

**Driver fatigue:** Available data are images from many driving videos of some group of people

Which of these are good **train/validation splits**?

▶ split the images randomly into two subsets, one for training and the other for validation ✗

▶ split the videos randomly into two subsets... ?  → what about repeated drivers?

▶ split the drivers randomly into two subsets... ?  → repeated car models?

Random splits are hard to make reproducible. How about deterministic ones?

▶ split the videos by time stamp...

**Driver fatigue:** Available data are images from many driving videos of some group of people

Which of these are good **train/validation splits**?

- ▶ split the images randomly into two subsets, one for training and the other for validation ✗
- ▶ split the videos randomly into two subsets... ？ → what about repeated drivers?
- ▶ split the drivers randomly into two subsets... ？ → repeated car models?

Random splits are hard to make reproducible. How about deterministic ones?
- ▶ split the videos by time stamp... ✗

## Quiz: train/validation sets

**Driver fatigue:** Available data are images from many driving videos of some group of people

Which of these are good **train/validation splits**?

▶ split the images randomly into two subsets, one for training and the other for validation ✗

▶ split the videos randomly into two subsets. . .   ?   → what about repeated drivers?

▶ split the drivers randomly into two subsets. . .   ?   → repeated car models?

Random splits are hard to make reproducible. How about deterministic ones?

▶ split the videos by time stamp. . . ✗

▶ split the drivers by gender. . .

**Driver fatigue:** Available data are images from many driving videos of some group of people

Which of these are good **train/validation splits**?

- ▶ split the images randomly into two subsets, one for training and the other for validation ✗
- ▶ split the videos randomly into two subsets... ? → what about repeated drivers?
- ▶ split the drivers randomly into two subsets... ? → repeated car models?

Random splits are hard to make reproducible. How about deterministic ones?
- ▶ split the videos by time stamp... ✗
- ▶ split the drivers by gender... ✗

**Driver fatigue:** Available data are images from many driving videos of some group of people

Which of these are good **train/validation splits**?

- ▶ split the images randomly into two subsets, one for training and the other for validation ✗
- ▶ split the videos randomly into two subsets... ? → what about repeated drivers?
- ▶ split the drivers randomly into two subsets... ? → repeated car models?

Random splits are hard to make reproducible. How about deterministic ones?

- ▶ split the videos by time stamp... ✗
- ▶ split the drivers by gender... ✗
- ▶ split the drivers by last name (A-L vs M-Z)

**Driver fatigue:** Available data are images from many driving videos of some group of people

Which of these are good **train/validation splits**?

- ▶ split the images randomly into two subsets, one for training and the other for validation ✗
- ▶ split the videos randomly into two subsets... ? → what about repeated drivers?
- ▶ split the drivers randomly into two subsets... ? → repeated car models?

Random splits are hard to make reproducible. How about deterministic ones?

- ▶ split the videos by time stamp... ✗
- ▶ split the drivers by gender... ✗
- ▶ split the drivers by last name (A–L vs M–Z) ? → can introduce bias

**Driver fatigue:** Available data are images from many driving videos of some group of people

Which of these are good **train/validation splits**?

- ▶ split the images randomly into two subsets, one for training and the other for validation ✗

- ▶ split the videos randomly into two subsets... ❓ → what about repeated drivers?

- ▶ split the drivers randomly into two subsets... ❓ → repeated car models?

Random splits are hard to make reproducible. How about deterministic ones?

- ▶ split the videos by time stamp... ✗

- ▶ split the drivers by gender... ✗

- ▶ split the drivers by last name (A-L vs M-Z) ❓ → can introduce bias

**Summary: It's not always obvious how to form a validation set (especially for time series/videos).**

## Model selection with a validation set

Given: training set $S = \{(x_1, y_1), \ldots, (x_n, y_n)\}$, possible hyperparameters $A = \{\eta_1, \ldots, \eta_m\}$ (including, e.g., which model class to use).

▶ Split available training data in to disjoint *real train* and *validation* set

$$S = S_{\text{trn}} \;\dot\cup\; S_{\text{val}}$$

▶ for all $\eta \in A$:

▶    $f^\eta \leftarrow$ train a model with hyperparameters $\eta$ using data $S_{\text{trn}}$

▶    $E_{\text{val}}^\eta \leftarrow$ evaluate model $f^\eta$ on set $S_{\text{val}}$

▶ $\eta^* = \operatorname{argmin}_{\eta \in A} E_{\text{val}}^\eta$      (select most promising hyperparameters)

▶ optionally: retrain model with hyperparameters $\eta^*$ on complete $S$

degree 1 fit, $\hat{\mathcal{R}} = 11.16$

true signal $\mathcal{R} = 14.33$

training points

validation points $\hat{\mathcal{R}}_{val} = 14.36$

degree 6 fit, $\hat{\mathcal{R}} = 0.03$
true signal $\mathcal{R} = 0.31$
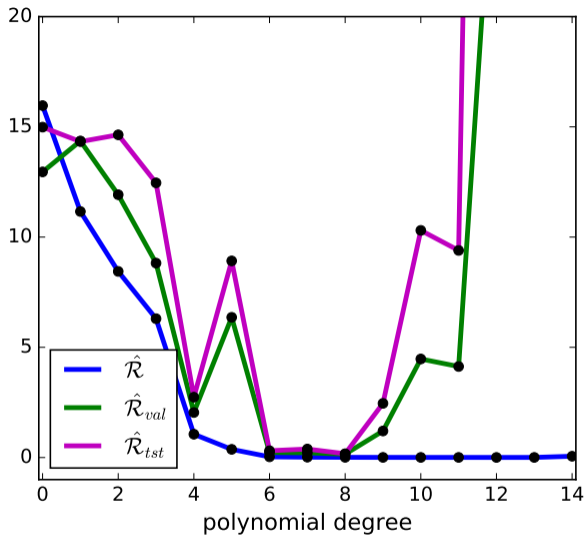training points
validation points $\hat{\mathcal{R}}_{val} = 0.23$

# Illustration: learning polynomials (of different degrees)

| degree | $\hat{\mathcal{R}}$ | $\hat{\mathcal{R}}_{val}$ | $\mathcal{R}$ |
|--------|------|-------|--------|
| 0  | 15.96 | 12.96   | 14.99   |
| 1  | 11.16 | 14.36   | 14.33   |
| 2  | 8.44  | 11.92   | 14.64   |
| 3  | 6.30  | 8.82    | 12.46   |
| 4  | 1.06  | 2.04    | 2.74    |
| 5  | 0.37  | 6.36    | 8.92    |
| 6  | 0.03  | 0.23    | 0.31    |
| 7  | 0.02  | 0.22    | 0.39    |
| 8  | 0.01  | 0.15    | 0.17    |
| 9  | 0.01  | 1.20    | 2.46    |
| 10 | 0.01  | 4.48    | 10.31   |
| 11 | 0.01  | 2.91    | 6.06    |
| 12 | 0.00  | 30.34   | 104.11  |
| 13 | 0.00  | 40.87   | 142.73  |
| 14 | 0.00  | 1622.37 | 8494.42 |

| degree | $\hat{\mathcal{R}}$ | $\hat{\mathcal{R}}_{val}$ | $\mathcal{R}$ |
|--------|------|---------|---------|
| 0 | 15.96 | 12.96 | 14.99 |
| 1 | 11.16 | 14.36 | 14.33 |
| 2 | 8.44 | 11.92 | 14.64 |
| 3 | 6.30 | 8.82 | 12.46 |
| 4 | 1.06 | 2.04 | 2.74 |
| 5 | 0.37 | 6.36 | 8.92 |
| 6 | 0.03 | 0.23 | 0.31 |
| 7 | 0.02 | 0.22 | 0.39 |
| **8** | **0.01** | **0.15** | **0.17** |
| 9 | 0.01 | 1.20 | 2.46 |
| 10 | 0.01 | 4.48 | 10.31 |
| 11 | 0.01 | 2.91 | 6.06 |
| 12 | 0.00 | 30.34 | 104.11 |
| 13 | 0.00 | 40.87 | 142.73 |
| 14 | 0.00 | 1622.37 | 8494.42 |

## From infinite to finite hyperparameter set

Typically, finitely many choices for model classes, optimizers
- ▶ we can/have to try them all

For hyperparameters typically too many or even infinitely many choices:
- ▶ neural network depth and widths
- ▶ learning rate(s)
- ▶ regularization strength
- ▶ number of training iterations

Subsample or discretize in a reasonable way, e.g.
- ▶ regularization: exponentially, e.g. $\lambda \in \{2^{-20}, 2^{-19}, \ldots, 2^{20}\}$
- ▶ learning rate: exponential scale in known range, e.g. $\eta \in \{10^{-6}, 3 \cdot 10^{-6}, \ldots, 10^{-1}\}$
- ▶ iterations/epochs: linearly, e.g. $T \in \{1, 5, 10, 50, 100, \ldots\}$ (or adaptive using $\mathcal{R}_{\mathsf{val}}$)

If model selection picks value at boundary of range, choose a larger range.

Selecting multiple hyperparameters: $\eta_1 \in A_1, \eta_2 \in A_2, \ldots, \eta_J \in A_J$

Selecting multiple hyperparameters: $\eta_1 \in A_1, \eta_2 \in A_2, \ldots, \eta_J \in A_J$

Grid search (good, but rarely practical for more than $J > 2$)

- $\eta_{\text{total}} = (\eta_1, \ldots, \eta_K)$, $A_{\text{total}} = A_1 \times A_2 \times \cdots \times A_J$
- try all $|A_1| \times |A_2| \times \cdots \times |A_J|$ combinations

Selecting multiple hyperparameters: $\eta_1 \in A_1, \eta_2 \in A_2, \ldots, \eta_J \in A_J$

Grid search (good, but rarely practical for more than $J > 2$)

- $\eta_{\text{total}} = (\eta_1, \ldots, \eta_K)$, $A_{\text{total}} = A_1 \times A_2 \times \cdots \times A_J$
- try all $|A_1| \times |A_2| \times \cdots \times |A_J|$ combinations

Greedy search (sometimes practical, but suboptimal)

- for $j = 1, \ldots, J$:
- $\eta_j^* \leftarrow$ pick $\eta_j \in A_j$, with $\eta_k$ for $k > j$ fixed at a reasonable default

Problem: what are 'reasonable defaults'?

Selecting multiple hyperparameters: $\eta_1 \in A_1, \eta_2 \in A_2, \ldots, \eta_J \in A_J$

## Grid search (good, but rarely practical for more than $J > 2$)

- $\eta_{\text{total}} = (\eta_1, \ldots, \eta_K)$, $A_{\text{total}} = A_1 \times A_2 \times \cdots \times A_J$
- try all $|A_1| \times |A_2| \times \cdots \times |A_J|$ combinations

## Greedy search (sometimes practical, but suboptimal)

- for $j = 1, \ldots, J$:
- $\eta_j^* \leftarrow$ pick $\eta_j \in A_j$, with $\eta_k$ for $k > j$ fixed at a reasonable default

Problem: what are 'reasonable defaults'?

## Trust in a higher authority (avoid this!)

- set hyperparameters to values from the literature, e.g. $\lambda = 1$

Problem: this can fail horribly if the situation isn't completely identical.

**Be sceptical of models with many ($> 2$) hyperparameters!**

# Common traps: don't fall for these!

## Trap 2: domain adaptation

Imagine you develop a method for *domain adaptation (DA)*:

▶ given: labeled data from some data distribution ("source")

▶ goal: create a classifier that works for a different data distribution ("target"), from which only unlabeled data is given

**How to do model-selection?**

## Trap 2: domain adaptation

Imagine you develop a method for *domain adaptation (DA)*:

- ▶ given: labeled data from some data distribution ("source")
- ▶ goal: create a classifier that works for a different data distribution ("target"), from which only unlabeled data is given

**How to do model-selection?**

**Reverse validation:**

- ▶ train DA classifier on labeled source, use it to predict labels for unlabeled target data
- ▶ train DA classifier on (now labeled) target data, predict on source
- ▶ compare predicted labels and ground truth labels on source

Tedious and not perfect, but better than nothing...

---

[E. Zhong, W. Fan, Q. Yang, O. Verscheure, J. Ren. "Cross validation framework to choose amongst models and datasets for transfer learning." ECML/PKDD 2010], [L. Bruzzone, M. Marconcini, "Domain adaptation problems: A DASVM classification technique and a circular validation strategy", TPAMI 2010].

Trap 3: zero-shot learning

Imagine: at prediction time, *classes will occur that are not in the training set*, e.g.

- ▶ training data: images of 40 animal classes
- ▶ prediction time: images of 10 other animals classes

**How to do model-selection?**

## Trap 3: zero-shot learning

Imagine: at prediction time, *classes will occur that are not in the training set*, e.g.

- ▶ training data: images of 40 animal classes
- ▶ prediction time: images of 10 other animals classes

**How to do model-selection?**

**Simulate prediction situation:**

- ▶ split training data into actual train and val data with disjoint classes (e.g. 32-8)
- ▶ train zero-shot classifier on actual train data to predict validation classes
- ▶ evaluate on val
- ▶ potential repeat multiplie times with different splits and average

Tedious to implement, and needs sufficiently many classes, but does the trick...

[Y. Xian, CHL, B. Schiele, Z. Akata, "Zero-shot learning-A comprehensive evaluation of the good, the bad and the ugly", TPAMI 2018]

Trap 4: outlier detection / out-of-distribution prediction

Imagine you want a system that knows what data it is useful for

▶ goal: classification of 10 object classes

▶ task: if other data occurs, refuse to classify

**How to do model-selection?**

Trap 4: outlier detection / out-of-distribution prediction

Imagine you want a system that knows what data it is useful for

- ▶ goal: classification of 10 object classes
- ▶ task: if other data occurs, refuse to classify

**How to do model-selection?**

**Unclear.**

- ▶ option 1: make up "out-of-distribution" data, e.g. random noise
  Problem: system might only detect noise, not other out-of-distribution data
- ▶ option 2: leave out some classes, use those as 'outliers' (like zero-shot)
  Problem: system might only detect new object classes, not e.g. pure noise
- ▶ option 3: use optimization to create "difficult" data, like adversarial examples
  Problem: very hard optimization problem, might still not work

**Same problem for actual evaluation, because no test set exists.**

## Trap 5: weakly-supervised learning

Imagine you learn from *weakly annotated data*, e.g.

▶ goal: object detection, i.e. predicting object bounding boxes

▶ given: training image with per-image labels

**How to do model-selection?**

Trap 5: weakly-supervised learning

Imagine you learn from *weakly annotated data*, e.g.

▶ goal: object detection, i.e. predicting object bounding boxes
▶ given: training image with per-image labels

**How to do model-selection?**

**I have no clue.**

▶ option 1: choose model and hyperparameters heuristically, hope for the best
▶ option 2: choose model and hyperparameters on another dataset that has bounding box annotation, hope for the best
▶ option 3: admit that you will need some images with bounding box annotation for model selection → but: why not use (part of) this for training?

If you can come up with a better way, please let me know.

## Summary: Machine Learning for Computer Vision

Step 1) Decide what exactly you want: the computer can't read your thoughts (yet)

▶ **inputs** $x$, **outputs** $y$, **loss** $\ell$, **model class** $f_\theta$ (with hyperparameters)

Step 2) Collect and annotate data: a good model with bad data is useless

▶ **collect and annotate data**, ideally i.i.d. from the prediction-time data distribution

Step 3) Model training ← often repeatedly

▶ **model training** on a training data, **model selection** on validation subset

Step 4) Model evaluation

▶ **evaluate model** on test set (that is disjoint from training/validation set)

# The Big Picture

**What's the purpose of (academic) science/research?**

**What's the purpose of (academic) science/research?**

- getting rich and famous?

**What's the purpose of (academic) science/research?**

▶ getting rich and famous? ✗

**What's the purpose of (academic) science/research?**

▶ getting rich and famous? ✗

▶ having fun while being paid for it?

**What's the purpose of (academic) science/research?**

▶ getting rich and famous?   ✗

▶ having fun while being paid for it?   ?

**What's the purpose of (academic) science/research?**

▶ getting rich and famous?  ✗

▶ having fun while being paid for it?  ?

▶ pleasing your supervisor?

**What's the purpose of (academic) science/research?**

▶ getting rich and famous? ✗

▶ having fun while being paid for it? ?

▶ pleasing your supervisor? ?

**What's the purpose of (academic) science/research?**

- ▶ getting rich and famous? ✗
- ▶ having fun while being paid for it? ?
- ▶ pleasing your supervisor? ?
- ▶ advancing human knowledge and understanding?

**What's the purpose of (academic) science/research?**

▶ getting rich and famous?  ✗

▶ having fun while being paid for it?  ?

▶ pleasing your supervisor?  ?

▶ advancing human knowledge and understanding?  ✓

**Two often overlooked aspects:**

▶ You're doing research **for others**, not just for yourself.
  – do something good, and then tell others about it

▶ You have to **advance** something for this to make sense.
  – it's not actually about doing something good, it's about doing something <u>better</u>.

**What's the purpose of (academic) science/research?**

▶ getting rich and famous?  ✗

▶ having fun while being paid for it?  ?

▶ pleasing your supervisor?  ?

▶ advancing human knowledge and understanding?  ✓

**Two often overlooked aspects:**

▶ You're doing research **for others**, not just for yourself.
  – do something good, and then tell others about it ... and they have to believe you!

▶ You have to **advance** something for this to make sense.
  – it's not actually about doing something good, it's about doing something <u>better</u>.

**One of the cornerstones of science is that results must be reproducible:**

Any person knowledgeable in the field should be able to repeat your experiments and get the same results.

### Reproducibility

▶ reproducing the experiments requires access to the data
  – use public data and provide information where you got it from
  – if you have to use your own data, make it publicly available

▶ reproducing the experiments requires knowledge of all components
  – clearly describe all components used, not just the new ones
  – list all implementation details (preferably release the source code)
  – list all (hyper)parameter values and how they were obtained

**Your goal should not be to just to get papers accepted, but to create new knowledge.**

Be your own method's harshest critic: you know best what's going on under the hood and what could have gone wrong in the process.

### Scientific scrutiny

- ▶ are you solving a relevant problem?
- ▶ is the data representative for the actual problem?
- ▶ are the results explainable just by random chance/noise?
- ▶ does the system make use of artifacts in the data?
- ▶ can you rule out bugs in the code or data leakage?
- ▶ if you compare to baselines, is the comparison *fair*?

Just that you invested a lot of work into a project doesn't mean it deserves to be published. A project must be able to *fail*, otherwise it's not research.

## Scientific scrutiny

- ▶ are you solving a relevant problem?
- ▶ is the data representative for the actual problem?
- ▶ are the results explainable just by random chance/noise?
- ▶ does the system make use of artifacts in the data?
- ▶ can you rule out bugs in the code or data leakage?
- ▶ if you compare to baselines, are they truly comparable?

Machine learning can be used for many purposes. Not all are desirable.

DOI: 10.1037/pspa0000098 · Corpus ID: 1379347

## Deep Neural Networks Are More Accurate Than Humans at Detecting Sexual Orientation From Facial Images

Yilun Wang, M. Kosinski · Published 1 February 2018 · Computer Science · Journal of Personality and Social Psychology

We show that faces contain much more information about sexual orientation than can be perceived or interpreted by the human brain. We used deep neural networks to extract features from 35,326 facial images. These features were entered into a logistic regression aimed at classifying sexual orientation. Given a single facial image, a classifier could correctly distinguish between gay and heterosexual men in 81% of cases, and in 71% of cases for women. Human judges achieved much lower accuracy: 61… Expand

## Scientific scrutiny

- ▶ are you solving a relevant problem?
- ▶ is the data representative for the actual problem?
- ▶ are the results explainable just by random chance/noise?
- ▶ does the system make use of artifacts in the data?
- ▶ can you rule out bugs in the code or data leakage?
- ▶ if you compare to baselines, are they truly comparable?

Example task: emotion recognition for autonomous driving
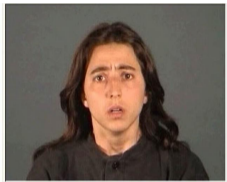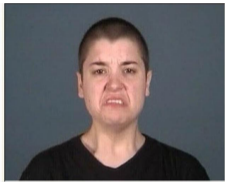- ▶ e.g. $99\%$ accuracy on DAFEX database for facial expressions

Example task: emotion recognition for autonomous driving
▶ e.g. $99\%$ accuracy on DAFEX database for facial expressions
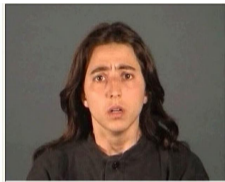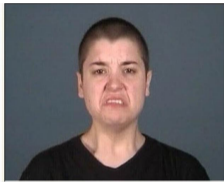


| Sadness | Disgust | Happiness | Surprise |

Example task: emotion recognition for autonomous driving

▶ e.g. $99\%$ accuracy on DAFEX database for facial expressions



**Sadness**  **Disgust**  **Happiness**  **Surprise**

▶ frontal faces at fixed distance. Will that be the case in the car?

Example task: emotion recognition for autonomous driving
- ▶ e.g. $99\%$ accuracy on DAFEX database for facial expressions



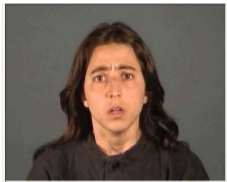**Sadness**          **Disgust**          **Happiness**          **Surprise**

- ▶ frontal faces at fixed distance. Will that be the case in the car?
- ▶ perfect lighting conditions. Will that be the case in the car?

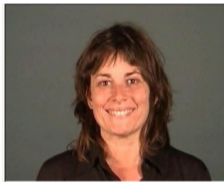Example task: emotion recognition for autonomous driving

▶ e.g. $99\%$ accuracy on DAFEX database for facial expressions



**Sadness**  **Disgust**  **Happiness**  **Surprise**

▶ frontal faces at fixed distance. Will that be the case in the car?
▶ perfect lighting conditions. Will that be the case in the car?
▶ homogeneous gray background. Will that be the case in the car?

Example task: emotion recognition for autonomous driving
- e.g. $99\%$ accuracy on DAFEX database for facial expressions



**Sadness**  **Disgust**  **Happiness**  **Surprise**

- frontal faces at fixed distance. Will that be the case in the car?
- perfect lighting conditions. Will that be the case in the car?
- homogeneous gray background. Will that be the case in the car?
- all faces are Italian actors. Will all drivers be Italian actors?

Example task: emotion recognition for autonomous driving

▶ e.g. $99\%$ accuracy on DAFEX database for facial expressions



**Sadness**     **Disgust**     **Happiness**     **Surprise**

▶ frontal faces at fixed distance. Will that be the case in the car?
▶ perfect lighting conditions. Will that be the case in the car?
▶ homogeneous gray background. Will that be the case in the car?
▶ all faces are Italian actors. Will all drivers be Italian actors?
▶ the emotions are not real, but (over)acted!

Example task: emotion recognition for autonomous driving
▶ e.g. $99\%$ accuracy on DAFEX database for facial expressions



**Sadness**  **Disgust**  **Happiness**  **Surprise**

▶ frontal faces at fixed distance. Will that be the case in the car?
▶ perfect lighting conditions. Will that be the case in the car?
▶ homogeneous gray background. Will that be the case in the car?
▶ all faces are Italian actors. Will all drivers be Italian actors?
▶ the emotions are not real, but (over)acted!

**Little reason to believe that the system will work in practice → don't overclaim!**

Images: https://i3.fbk.eu/resources/dafex-database-kinetic-facial-expressions

## Scientific scrutiny

- ▶ are you solving a relevant problem?
- ▶ is the data representative for the actual problem?
- ▶ are the results explainable just by random chance/noise?
- ▶ does the system make use of artifacts in the data?
- ▶ can you rule out bugs in the code or data leakage?
- ▶ if you compare to baselines, are they truly comparable?

## Repeated Experiments / Randomization

In the natural sciences, to show that an effect was not just random chance one uses repeated experiments.

In computer science, running the same code multiple times should give identical results, so on top we must use a form of explicit randomization.

▶ different random splits of the data, e.g. for training / evaluation
▶ different random initialization, e.g. for neural network training

Which of these make sense depends on the situation

▶ for benchmark datasets, test set is usually fixed
▶ for convex optimization methods, initialization plays no role, etc.

## Repeats

We can think of accuracy on a validation/test set also as repeated experiments:
- ▶ we apply a fixed model many times, each time to a single example
- ▶ outcomes: $r_1, r_2, \ldots, r_m$

Having done repeated experiments, one reports either
- ▶ the mean of outcomes and standard deviations

$$\bar{r} \pm \sigma \qquad \text{for} \qquad \bar{r} = \frac{1}{m} \sum_j r_t, \qquad \sigma = \sqrt{\frac{1}{m} \sum_j (r_j - \bar{r})^2}$$

*"If you try this once, where can you expect the result to lie?"*

or
- ▶ the mean and standard error of the mean

$$\bar{r} \pm \text{SE} \qquad \text{for} \qquad \bar{r} = \frac{1}{m} \sum_j r_t, \quad \text{SE} = \frac{\sigma}{\sqrt{m}}$$

*"If you try this many times, where can you expect the mean to lie?"*

Illustration as figure with error bars:



Illustration as table with error intervals:

| error rate | Chimpanzee | Giant panda |
|---|---|---|
| component 1 | $25.47 \pm 0.42$ | $21.02 \pm 0.58$ |
| component 2 | $23.94 \pm 0.32$ | $19.44 \pm 0.50$ |
| combination | $23.66 \pm 0.33$ | $19.23 \pm 0.52$ |

Quick (and dirty) check if results could be due to random chance:
do mean$\pm$error bars overlap?

# Box and Whisker Plot

More informative: **box and whisker plots** (short: box plot)



- red line: data median
- blue body: 25%–75% quantile
- whiskers: value range
- blue markers: outliers

sometimes additional symbols for mean or uncertainty of the median

We can think of accuracy on a test set as repeated experiments:
apply a fixed model many times, each time to a single example.

Which of these differences are real and which are random fluctuations?

## Quiz: Randomness in Repeated Experiments

We can think of accuracy on a test set as repeated experiments:
apply a fixed model many times, each time to a single example.

Which of these differences are real and which are random fluctuations?

- ▶ method A: 71% accuracy
- ▶ method B: 85% accuracy

We can think of accuracy on a test set as repeated experiments:
apply a fixed model many times, each time to a single example.

Which of these differences are real and which are random fluctuations?

- method A: 71% accuracy
- method B: 85% accuracy

unclear: 5/7 vs 6/7 or 710/1000 vs 850/1000?
always report test set size!

We can think of accuracy on a test set as repeated experiments:
apply a fixed model many times, each time to a single example.

Which of these differences are real and which are random fluctuations?

- method A: 71% accuracy
- method B: 85% accuracy

unclear: 5/7 vs 6/7 or 710/1000 vs 850/1000?
always report test set size!

- A: 40% accuracy on 10 examples
- B: 60% accuracy on 10 examples

We can think of accuracy on a test set as repeated experiments:
apply a fixed model many times, each time to a single example.

Which of these differences are real and which are random fluctuations?

- method A: 71% accuracy
- method B: 85% accuracy

unclear: 5/7 vs 6/7 or 710/1000 vs 850/1000?
always report test set size!

- A: 40% accuracy on 10 examples
- B: 60% accuracy on 10 examples

high chance that effect is random ($> 14\%$)

We can think of accuracy on a test set as repeated experiments:
apply a fixed model many times, each time to a single example.

Which of these differences are real and which are random fluctuations?

- method A: 71% accuracy
- method B: 85% accuracy

unclear: 5/7 vs 6/7 or 710/1000 vs 850/1000?
always report test set size!

- A: 40% accuracy on 10 examples
- B: 60% accuracy on 10 examples

high chance that effect is random ($> 14\%$)

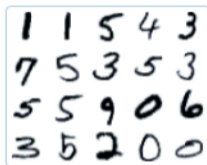- A: 40% accuracy on 1000 examples
- B: 60% accuracy on 1000 examples

We can think of accuracy on a test set as repeated experiments:
apply a fixed model many times, each time to a single example.

Which of these differences are real and which are random fluctuations?

- method A: 71% accuracy
- method B: 85% accuracy

unclear: 5/7 vs 6/7 or 710/1000 vs 850/1000?
always report test set size!

- A: 40% accuracy on 10 examples
- B: 60% accuracy on 10 examples

high chance that effect is random ($> 14\%$)

- A: 40% accuracy on 1000 examples
- B: 60% accuracy on 1000 examples

small chance that effect is random ($< 10^{-10}$)

We can think of accuracy on a test set as repeated experiments:
apply a fixed model many times, each time to a single example.

Which of these differences are real and which are random fluctuations?

- method A: 71% accuracy
- method B: 85% accuracy

unclear: 5/7 vs 6/7 or 710/1000 vs 850/1000?
always report test set size!

- A: 40% accuracy on 10 examples
- B: 60% accuracy on 10 examples

high chance that effect is random $(> 14\%)$

- A: 40% accuracy on 1000 examples
- B: 60% accuracy on 1000 examples

small chance that effect is random $(< 10^{-10})$

- A: 99.77% accuracy on 10000 ex.s
- B: 99.79% accuracy on 10000 ex.s

We can think of accuracy on a test set as repeated experiments:
apply a fixed model many times, each time to a single example.

Which of these differences are real and which are random fluctuations?

- method A: 71% accuracy
- method B: 85% accuracy

unclear: 5/7 vs 6/7 or 710/1000 vs 850/1000?
always report test set size!

- A: 40% accuracy on 10 examples
- B: 60% accuracy on 10 examples

high chance that effect is random ($> 14\%$)

- A: 40% accuracy on 1000 examples
- B: 60% accuracy on 1000 examples

small chance that effect is random ($< 10^{-10}$)

- A: 99.77% accuracy on 10000 ex.s
- B: 99.79% accuracy on 10000 ex.s

high chance that effect is random ($> 20\%$)

# MNIST

who is the best in MNIST ?



## MNIST 50 results collected

Units: error %

> Classify handwriten digits. Some additional results are available on the original dataset page.

| Result | Method | Venue | Details |
|--------|--------|-------|---------|
| 0.21% | Regularization of Neural Networks using DropConnect 📄 | ICML 2013 | |
| 0.23% | Multi-column Deep Neural Networks for Image Classification 📄 | CVPR 2012 | |
| 0.23% | APAC: Augmented PAttern Classification with Neural Networks 📄 | arXiv 2015 | |
| 0.24% | Batch-normalized Maxout Network in Network 📄 | arXiv 2015 | Details |

We can think of accuracy on a test set as repeated experiments:
apply a fixed model many times, each time to a single example.

Which of these differences are real and which are random fluctuations?

- method A: 71% accuracy
- method B: 85% accuracy

unclear. 5/7 vs 6/7 or 710/1000 vs 850/1000?

- A: 40% accuracy on 10 examples
- B: 60% accuracy on 10 examples

high chance that effect is random ($> 14\%$)

- A: 40% accuracy on 1000 examples
- B: 60% accuracy on 1000 examples

small chance that effect is random ($< 10^{-10}$)

- A: 99.77% accuracy on 10000 ex.s
- B: 99.79% accuracy on 10000 ex.s

high chance that effect is random ($> 20\%$)

Actual procedure: **statistical significance test!**

# Generalizability

In the natural sciences, to show that an effect is generally relevant one shows it in multiple model organisms.

In machine learning, this corresponds to multiple datasets and/or multiple models/architectures.

To show that a newly proposed model is better than earlier ones:

▶ show its superiority on multiple datasets

▶ ideally, show it for multiple different tasks (classification, segmentation, detection, . . . )

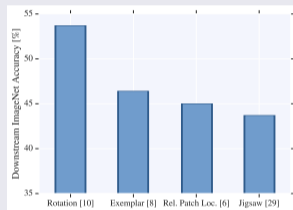To show that a new method or a new component is beneficial:

▶ show its superiority for multiple models

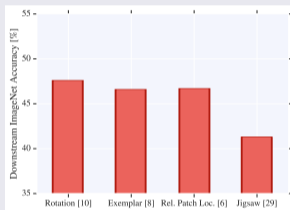▶ ideally, show it for multiple different tasks (classification, segmentation, detection, . . . )

In the natural sciences, to show that an effect is generally relevant one shows it in multiple model organisms.

In machine learning, this corresponds to multiple datasets and/or multiple models/architectures.

To show that a newly proposed model is better than earlier ones:

▶ show its superiority on multiple datasets

▶ ideally, show it for multiple different tasks (classification, segmentation, detection, . . . )

To show that a new method or a new component is beneficial:

▶ show its superiority for multiple models

▶ ideally, show it for multiple different tasks (classification, segmentation, detection, . . . )

**Datasets are proxies for real-world situations. Nobody should actually care about the results on the dataset.**

# Generalizability matters!

Observations: comparing on a single model architecture and/or dataset can be misleading.
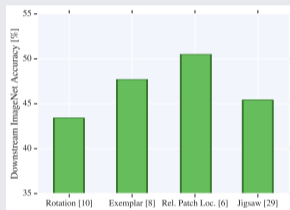
## Example: Unsupervised Pretraining for Image Classification



ResNet50        ResNet50v1        ResNet50v2

Method [10] looks great for original ResNet50, good for ResNet50v1, bad for ResNet50v2.

Designing convincing experiments is harder than it seems.

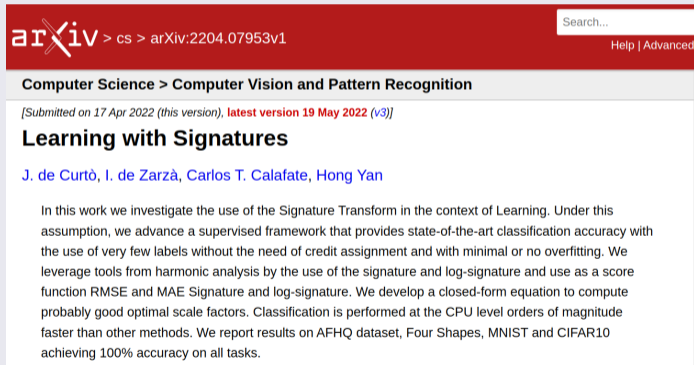Source: [A. Kolesnikov, X. Zhai, L. Beyer: "Revisiting Self-Supervised Visual Representation Learning", CVPR2019]

## Scientific scrutiny

► are you solving a relevant problem?
► is the data representative for the actual problem?
► are the results explainable just by random chance/noise?
► does the system make use of artifacts in the data?
► can you rule out bugs in the code or data leakage?
► if you compare to baselines, are they truly comparable?

🐎 **Springer** Open                    Search 🔍    Menu ▾

Journal of Big Data

Research | Open Access | Published: 07 January 2020

# Criminal tendency detection from facial images and the gender bias effect

Mahdi Hashemi ✉ & Margeret Hall

Paper reports 97% accuracy on 10.000 images (10-fold cross-validation).

Reported effect was explainable purely by a bias in the dataset collection.



criminals: *mugshots*          non-criminals: *normal photos from the web*

Note: original paper was withdrawn because the authors had not asked for ethics approval.

## Scientific scrutiny

▶ are you solving a relevant problem?
▶ is the data representative for the actual problem?
▶ are the results explainable just by random chance/noise?
▶ does the system make use of artifacts in the data?
▶ can you rule out bugs in the code or data leakage?
▶ if you compare to baselines, are they truly comparable?

**Always check your results for plausibility.**

Sometimes, good results are the effect of bugs or (test) data leakage.

If it looks too good to be true, it's probably not true.

**Computer Science > Computer Vision and Pattern Recognition**

[Submitted on 17 Apr 2022 (this version), latest version 19 May 2022 (v3)]
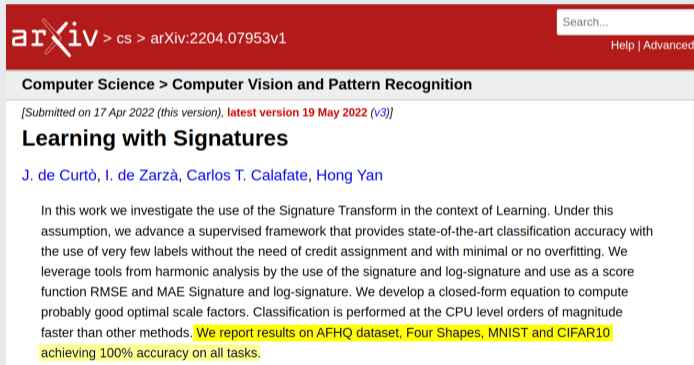
### Learning with Signatures

J. de Curtò, I. de Zarzà, Carlos T. Calafate, Hong Yan

In this work we investigate the use of the Signature Transform in the context of Learning. Under this assumption, we advance a supervised framework that provides state-of-the-art classification accuracy with the use of very few labels without the need of credit assignment and with minimal or no overfitting. We leverage tools from harmonic analysis by the use of the signature and log-signature and use as a score function RMSE and MAE Signature and log-signature. We develop a closed-form equation to compute probably good optimal scale factors. Classification is performed at the CPU level orders of magnitude faster than other methods. We report results on AFHQ dataset, Four Shapes, MNIST and CIFAR10 achieving 100% accuracy on all tasks.

**Always check your results for plausibility.**

Sometimes, good results are the effect of bugs or (test) data leakage.

If it looks too good to be true, it's probably not true.



arXiv > cs > arXiv:2204.07953v1

Search...

Help | Advanced

**Computer Science > Computer Vision and Pattern Recognition**

[Submitted on 17 Apr 2022 (this version), latest version 19 May 2022 (v3)]

## Learning with Signatures

J. de Curtò, I. de Zarzà, Carlos T. Calafate, Hong Yan

In this work we investigate the use of the Signature Transform in the context of Learning. Under this assumption, we advance a supervised framework that provides state-of-the-art classification accuracy with the use of very few labels without the need of credit assignment and with minimal or no overfitting. We leverage tools from harmonic analysis by the use of the signature and log-signature and use as a score function RMSE and MAE Signature and log-signature. We develop a closed-form equation to compute probably good optimal scale factors. Classification is performed at the CPU level orders of magnitude faster than other methods. We report results on AFHQ dataset, Four Shapes, MNIST and CIFAR10 achieving 100% accuracy on all tasks.

**Always check your results for plausibility.**
Sometimes, good results are the effect of bugs or (test) data leakage.

If it looks too good to be true, it's probably not true.



arXiv > cs > arXiv:2204.07953v1

Search...
Help | Advanced

**Computer Science > Computer Vision and Pattern Recognition**

[Submitted on 17 Apr 2022 (this version), latest version 19 May 2022 (v3)]

## Learning with Signatures

J. de Curtò, I. de Zarzà, Carlos T. Calafate, Hong Yan

In this work we investigate the use of the Signature Transform in the context of Learning. Under this assumption, we advance a supervised framework that provides state-of-the-art classification accuracy with the use of very few labels without the need of credit assignment and with minimal or no overfitting. We leverage tools from harmonic analysis by the use of the signature and log-signature and use as a score function RMSE and MAE Signature and log-signature. We develop a closed-form equation to compute probably good optimal scale factors. Classification is performed at the CPU level orders of magnitude faster than other methods. We report results on AFHQ dataset, Four Shapes, MNIST and CIFAR10 achieving 100% accuracy on all tasks.

cat?

cat?

**CIFAR test set has annotation errors, 100% should be impossible.**

Image: https://labelerrors.com/

Story: see https://www.reddit.com/r/MachineLearning/comments/u7ouxh/r_authors_claim_to_have_solved_mnist_and_cifar/

## Scientific scrutiny

- ▶ are you solving a relevant problem?
- ▶ is the data representative for the actual problem?
- ▶ are the results explainable just by random chance/noise?
- ▶ does the system make use of artifacts in the data?
- ▶ can you rule out bugs in the code or data leakage?
- ▶ if you compare to baselines, are they truly comparable?

Not always easy to answer:

## Do the methods have solve the same problem?

- ▶ previous work: object classification
- ▶ proposed work: object detection

Not always easy to answer:

### Do the methods have solve the same problem?

- ▶ previous work: object classification
- ▶ proposed work: object detection

### Do they evaluate their performance on the same data?

- ▶ previous work: PASCAL VOC 2010
- ▶ proposed work: PASCAL VOC 2012

## If you compare to baselines, are they truly comparable?

Not always easy to answer:

### Do the methods have solve the same problem?

- ▶ previous work: object classification
- ▶ proposed work: object detection

### Do they evaluate their performance on the same data?

- ▶ previous work: PASCAL VOC 2010
- ▶ proposed work: PASCAL VOC 2012

### Do they have access to the same data for training?

- ▶ previous work: ImageNet
- ▶ proposed work: ImageNet, pretrained on JFT-300M

# Do they have access to the same annotation?

- previous work: MSCOCO with bounding-box annotation
- proposed work: MSCOCO with segmentation annotation

## Do they have access to the same annotation?

- ▶ previous work: MSCOCO with bounding-box annotation
- ▶ proposed work: MSCOCO with segmentation annotation

## Do they use the same data augmentation?

- ▶ previous work: CIFAR-10 images
- ▶ proposed work: CIFAR-10 images with translation/rotation/crop augmentation

## Do they have access to the same annotation?

▶ previous work: MSCOCO with bounding-box annotation

▶ proposed work: MSCOCO with segmentation annotation

## Do they use the same data augmentation?

▶ previous work: CIFAR-10 images

▶ proposed work: CIFAR-10 images with translation/rotation/crop augmentation

## for runtime comparisons: Do you use comparable hardware?

▶ previous work: laptop with 3.4 GHz quad-core CPU

▶ proposed work: workstation with 8 A100 GPUs

## Do they have access to the same annotation?

▶ previous work: MSCOCO with bounding-box annotation
▶ proposed work: MSCOCO with segmentation annotation

## Do they use the same data augmentation?

▶ previous work: CIFAR-10 images
▶ proposed work: CIFAR-10 images with translation/rotation/crop augmentation

## for runtime comparisons: Do you use comparable hardware?

▶ previous work: laptop with 3.4 GHz quad-core CPU
▶ proposed work: workstation with 8 A100 GPUs

## for runtime comparisons: Do you use comparable languages/implementations?

▶ previous work: method implemented in pytorch
▶ proposed work: method implemented in tensorflow

## Is/was model selection for the baselines done thoroughly?

- ▶ baseline: hyperparameter set arbitrarily ad hoc
- ▶ proposed method: exhaustive hyperparameter selection

Now, is the new method better, or just the baseline worse than it could be?

## Is/was model selection for the baselines done thoroughly?

► baseline: hyperparameter set arbitrarily ad hoc

► proposed method: exhaustive hyperparameter selection

Now, is the new method better, or just the baseline worse than it could be?

Positive exception:

Nowozin, Bakır. "Discriminative Subsequence Mining for Action Classification", ICCV 2007

| Method | KTH accuracy |
|---|---|
| Niebles et al. [11], LOO, pLSA | 81.50 |
| Dollár et al. [4], LOO, SVM RBF | 80.66 |
| Schuldt et al. [17], splits, SVM match | 71.71 |
| Ke et al. [7], splits, forward feat.-sel. | 62.94 |
| | |
| Subsequence Boosting, $B = 12$, splits | 84.72 |

## Is/was model selection for the baselines done thoroughly?

- ▶ baseline: hyperparameter set arbitrarily ad hoc
- ▶ proposed method: exhaustive hyperparameter selection

Now, is the new method better, or just the baseline worse than it could be?

Positive exception:

Nowozin, Bakır. "Discriminative Subsequence Mining for Action Classification", ICCV 2007

| Method | KTH accuracy |
|---|---|
| Niebles et al. [11], LOO, pLSA | 81.50 |
| Dollár et al. [4], LOO, SVM RBF | 80.66 |
| Schuldt et al. [17], splits, SVM match | 71.71 |
| Ke et al. [7], splits, forward feat.-sel. | 62.94 |
| baseline SVM linear bin, 1-vs-1 | 83.33 |
| baseline SVM RBF bin, 1-vs-1 | 85.19 |
| baseline SVM $\chi^2$ bin, 1-vs-1 | 87.04 |
| Subsequence Boosting, $B = 12$, splits | 84.72 |

## Take home messages

Don't use machine learning unless you have good reasons.

Keep in mind: all data is random and all numbers are uncertain.

Learning is about finding a model that works on future data.

Proper model-selection is hard. Avoid free (hyper)parameters!

A test set is for evaluating a single final ultimate model, nothing else.

You don't do research for yourself, but for others (users, readers, society, . . . ).

Results must be convincing: honest, reproducible, not overclaimed.