



Machine Learning (69152)

Machine Learning Fundamentals



Happy to be here!



Content

DISCLAIMER: SLIDES WILL BE UPDATED DURING THE COURSE,
DOWNLOAD THE MOST UPDATED VERSION FROM MOODLE

- **Introduction**
- **K-nearest neighbours (as a simple example)**
- **Regression**
 - Linear regression
 - Non-linear regression
 - Overfitting, regularization, model selection
 - Robust regression
 - Ridge regression
- **Classification**
 - Logistic regression
 - SVM
- **Unsupervised methods**
 - PCA
- **Evaluation metrics**

Bibliography

- **Kevin P. Murphy, Probabilistic Machine Learning: An Introduction, The MIT Press, 2022** (available online <https://probml.github.io/pml-book/book1.html>)
- Marc Peter Deisenroth, A. Aldo Faisal, Cheng Soon Ong, Mathematics for Machine Learning, Cambridge University Press, 2020 (available online <https://mml-book.com/>)
- Christopher M. Bishop, Pattern Recognition and Machine Learning, Springer, 2006 (available online <https://www.microsoft.com/en-us/research/uploads/prod/2006/01/Bishop-Pattern-Recognition-and-Machine-Learning-2006.pdf>)
- Aurélien Géron, Hands-On Machine Learning with Scikit-Learn, Keras & TensorFlow, O'Reilly, 2019 (2nd edition)
- Online courses
 - Andrew Ng, Stanford CS229: Machine Learning,
<https://www.youtube.com/playlist?list=PLoROMvodv4rMiGQp3WXShMGgzqpfVfbU>
 - Ayush Singh and @freecodecamp, ML for beginners,
<https://www.youtube.com/watch?v=NWONeJKn6kc>
- Python resources
 - The Python tutorial: <https://docs.python.org/3/tutorial/>
 - numpy tutorial: <https://cs231n.github.io/python-numpy-tutorial/>
 - scikit-learn tutorials: <https://scikit-learn.org/stable/tutorial/index.html>

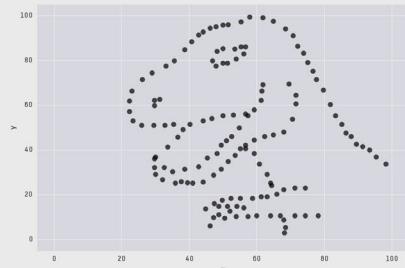
Motivation

- Machine learning has a huge presence (and growing)...
 - Recommenders, virtual assistants, content generation (intelligent virtual characters, automatic scene creation), spam filtering, search engines, social media, news classification, video surveillance, face recognition, fraud detection, traffic prediction, image retrieval, predictive maintenance, weather forecasting, autonomous driving, speech recognition, automatic translation, medical diagnosis...
- Starting by the fundamentals (and getting deep into math!) is critical for a structured understanding of the field.

$$\begin{aligned} &= \frac{n!}{k! \cdot (n-k)!} + \frac{n!}{(k+1)! \cdot (n-(k+1))!} \\ &= \frac{(k+1) \binom{n}{k} b_{n-k}^k}{(k+1) \cdot k! \cdot (n-k)!} + \frac{n! \cdot (n-k)}{(k+1)! \cdot (n-k)!} \\ &= \frac{(k+1)^k \binom{n-k}{k} + b \left(\sum_{k=0}^{n-k} \binom{n}{k} \cdot (n-k) \right)}{(k+1)! \cdot (n-k)!} \\ &= \frac{(k+1)^k \binom{n-k}{k} + \sum_{k=0}^{n-k} \binom{n}{k} a^k b^{n-k}}{(k+1) \cdot n! + \sum_{k=0}^{n-k} \binom{n}{k} a^k b^{n-k}} u_1^2 \\ &= \frac{(k+1)^k \binom{n-k}{k} + \sum_{k=0}^{n-k} \binom{n}{k} a^k b^{n-k}}{((k+1)^k + (n-k))^k} u_1^2 + P_1 + V_1 \\ &= \frac{(k+1)^k \binom{n-k}{k} + b^{n-k}}{((k+1)^k + (n-k))^k} u_1^2 + b^{n-k} \\ &= \frac{(k+1)^k \binom{n-k}{k} + b^{n-k}}{(n+1)^k + b^{n+1}} K = 1 - \sum_{n=1}^{\infty} \frac{1}{(2n-1)^5} \\ &= \frac{(n+1)!}{((n+1) \cdot p_0 \cdot (k+1)!)^2} = \frac{273.15}{\dots} \end{aligned}$$

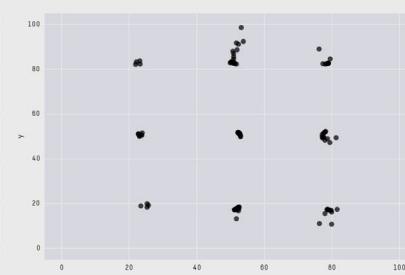
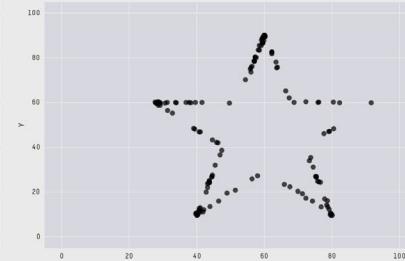
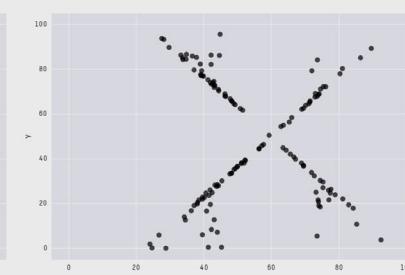
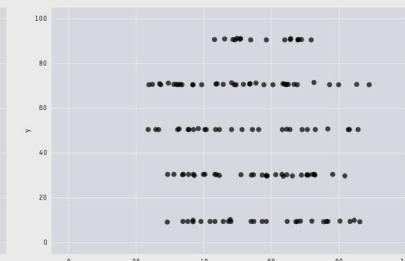
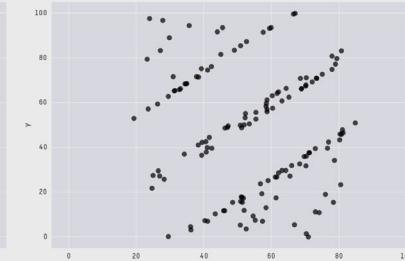
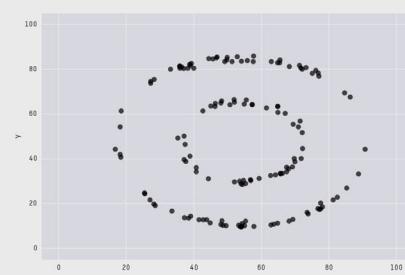
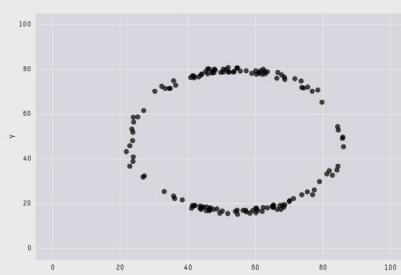
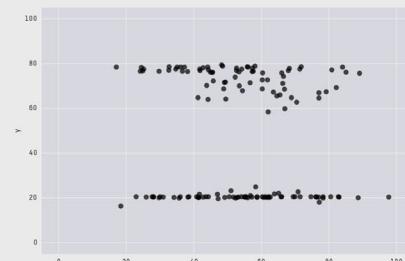
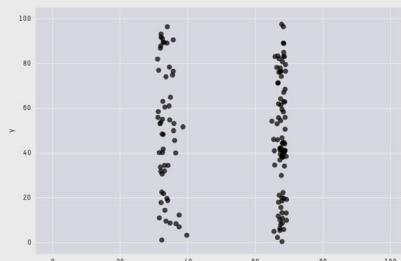
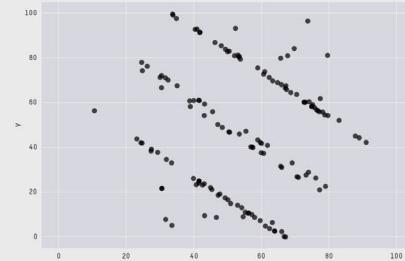
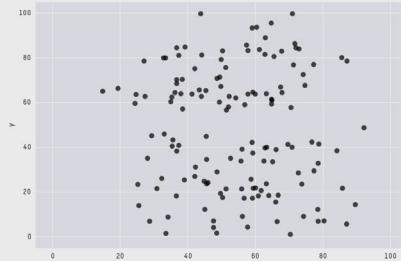
Harder than it seems!

Characterizing our data: what do we model? What is noise?



X Mean: 54.26
Y Mean: 47.83
X SD : 16.76
Y SD : 26.93
Corr. : -0.06

Alberto
Cairo, the
Datasaurus
dataset



Machine learning is NOT magic!

- Careful methodology, iterative design and data analysis can be motivated by the No Free Lunch theorem [Wolpert 96].
 - A model is a simplified version of the data
 - Assumptions (~hyperparameters) are needed (what should be accounted and what should be ignored in the data)
 - No model is guaranteed a priori to be better than others, it just should “fit” our data
 - Usual practice: Evaluate a few reasonable models making reasonable assumptions on your data.



Sometimes it seems magical though...

■ Example: Neural Radiance Fields

- Use neural networks (MLPs) for representing scenes and rendering novel views, costly to train and query

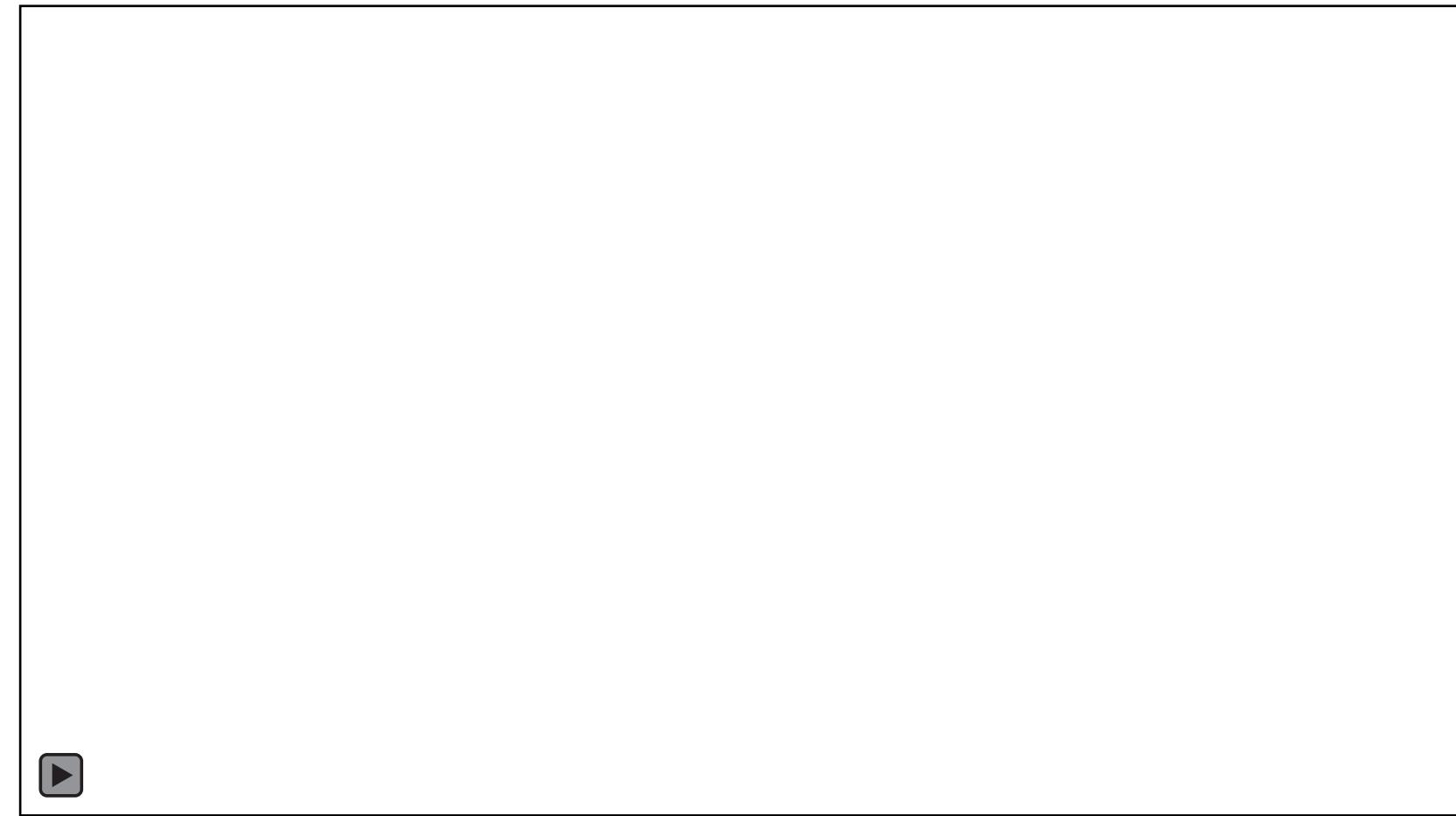


Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, Ren Ng, NeRF: Representing scenes as neural radiance fields for view synthesis, ECCV 2020
Thomas Müller and Alex Evans and Christoph Schied and Alexander Keller, Instant Neural Graphics Primitives with a Multiresolution Hash Encoding, ACM ToG, 2022

Sometimes it seems magical though...

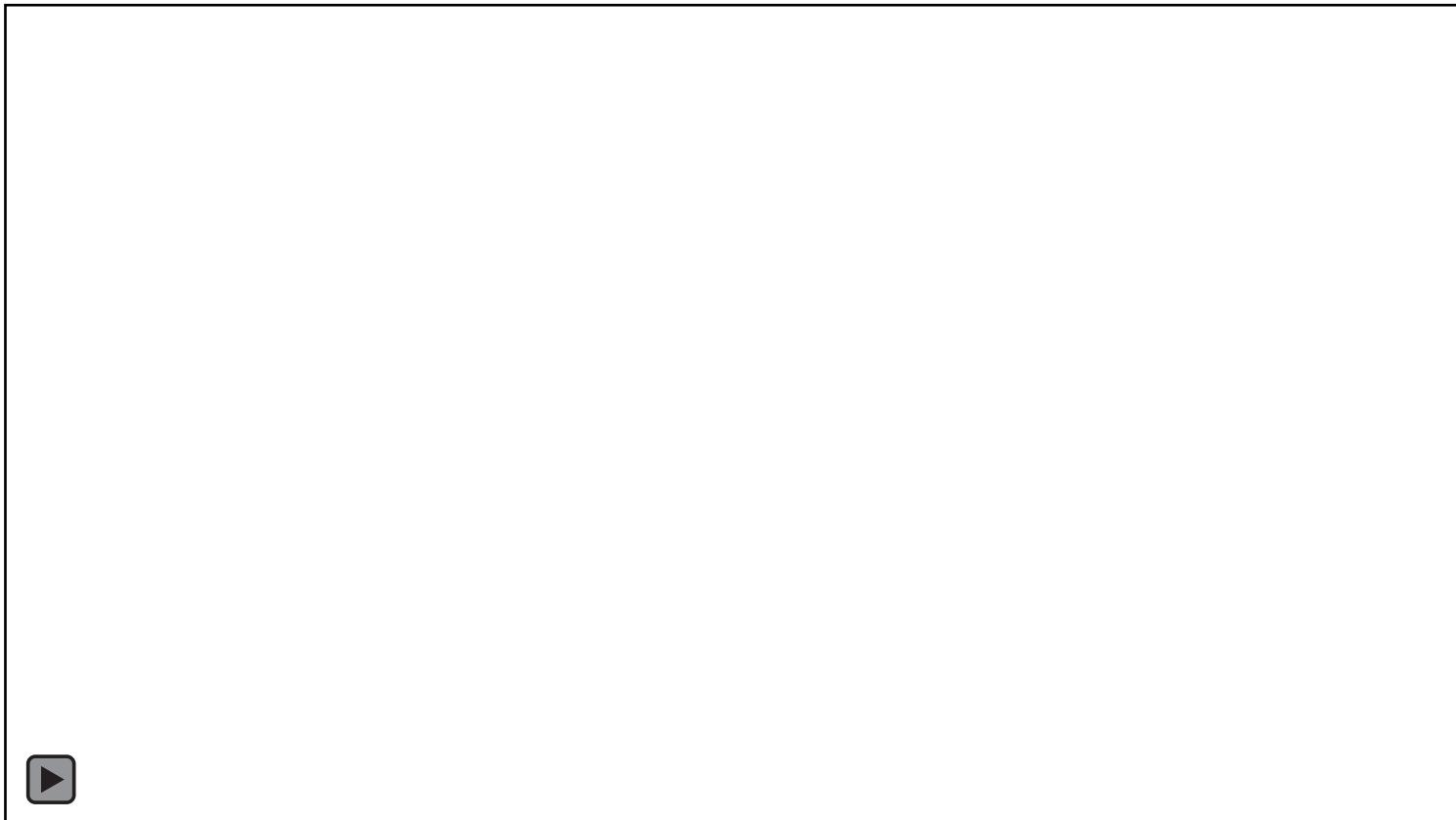
■ **Example:** Neural Radiance Fields

- Use neural networks (MLPs) for representing scenes and rendering novel views, costly to train and query



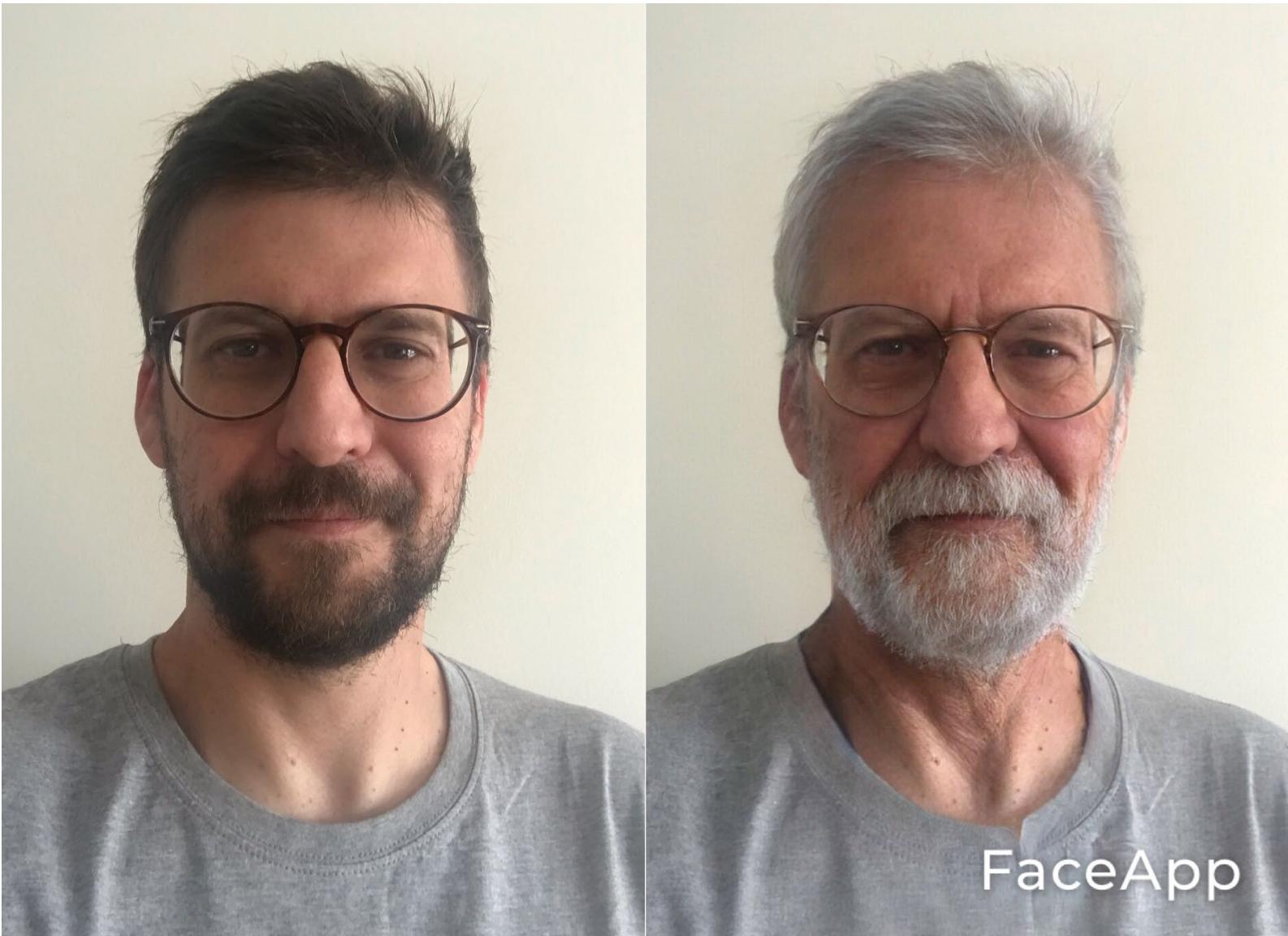
Shallow learning also works!

- Scene representation based on Gaussians in the 3D space
- Well implemented but standard point-based rendering
- Enables real-time capabilities at 1080p



Sometimes it seems magical though...

Just one cool example of sooo many...



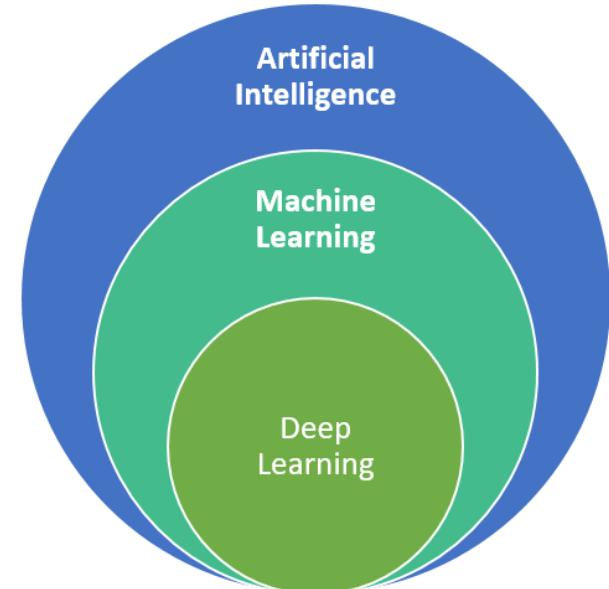
Definitions

■ Machine Learning:

- Field of study that gives computers the ability to learn without being explicitly programmed (Arthur Samuel 1959).
- A computer program is said to learn from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E. (Tom Mitchell, 1997)
- Success stories: Spam filtering, OCR, face detection, search engines...

■ Related terms:

- **Data mining:** The field of study that aims to extract information and structure from raw data.
- **Artificial Intelligence:** the theory and development of computer systems able to perform tasks normally requiring human intelligence, such as visual perception, speech recognition, decision-making, and translation between languages.
- **Knowledge Discovery, Pattern Recognition, ...**



Taxonomy

- **Training set:** Data samples that machine learning algorithms use to find the model parameters.
- **Test set:** Data used to assess generalization over samples that are unseen during training.
- **Supervised Learning:** The training set is composed of **labeled data** (data + desired output)
- **Unsupervised Learning:** The training data is **unlabeled**.
- **Regression:** Estimates a function that relates an input vector with a **continuous output**.
- **Classification:** Estimates a function that relates an input vector with a **discrete output**.
- **Parametric:** The model has a fixed number of parameters (by making assumptions on the data distribution)
- **Non-parametric:** The number of parameters grow with the data

Methodology (I de III)

1. Frame the Problem (Look at the Big Picture)

- Define the general high-level objective
- Frame the problem (supervised/unsupervised, regression/classification, etc.)
- Define metrics to evaluate performance (e.g., error, accuracy, false positives rate, etc.)
- What is a minimum/reasonable/maximum performance? Check scientific literature and existing similar systems.

2. Get the data

- What data is needed? How much? Which format?
- How/where do you get the data? Check legal aspects and sensitive information!
- Divide into training/validation/test sets. Never look at the test set!

Methodology (II de III)

3. Explore the data

- Visualize the data
- Study correlations between attributes

4. Prepare the data

- Work in a copy of the data, keep the original data intact
- Data cleaning: Fix/remove outliers, fill missing values / remove rows
- Feature selection: Drop uninformative attributes (garbage in – garbage out)
- Feature engineering: Transform the data into more descriptive/separable features (e.g., #murders vs. #murders/population)
- Feature scaling: Standardize or normalize features/data
- Shuffle the data

Methodology (III de III)

5. Shortlist promising models

- Sample subsets of the training sets for quick exploration of many techniques/models
- Compare models using cross-validation (mean and std over folds)
- Analyze errors

6. Fine-tune the system

- Fine-tune hyperparameters using cross-validation
- Try ensemble methods
- Once finished, evaluate the performance on the test set

7. Present the solution

- Report relevant aspects of the process, not only performance
- Failure cases might be relevant!
- Report assumptions, limitations, how the model could be improved
- Summarize (graphs, key findings)

K-nearest-neighbor (birds of a feather...)

- We have N sample pairs in our training dataset \mathcal{D}

$$\mathcal{D} = \{(x^{(1)}, y^{(1)}), \dots, (x^{(i)}, y^{(i)}), \dots, (x^{(N)}, y^{(N)})\}$$

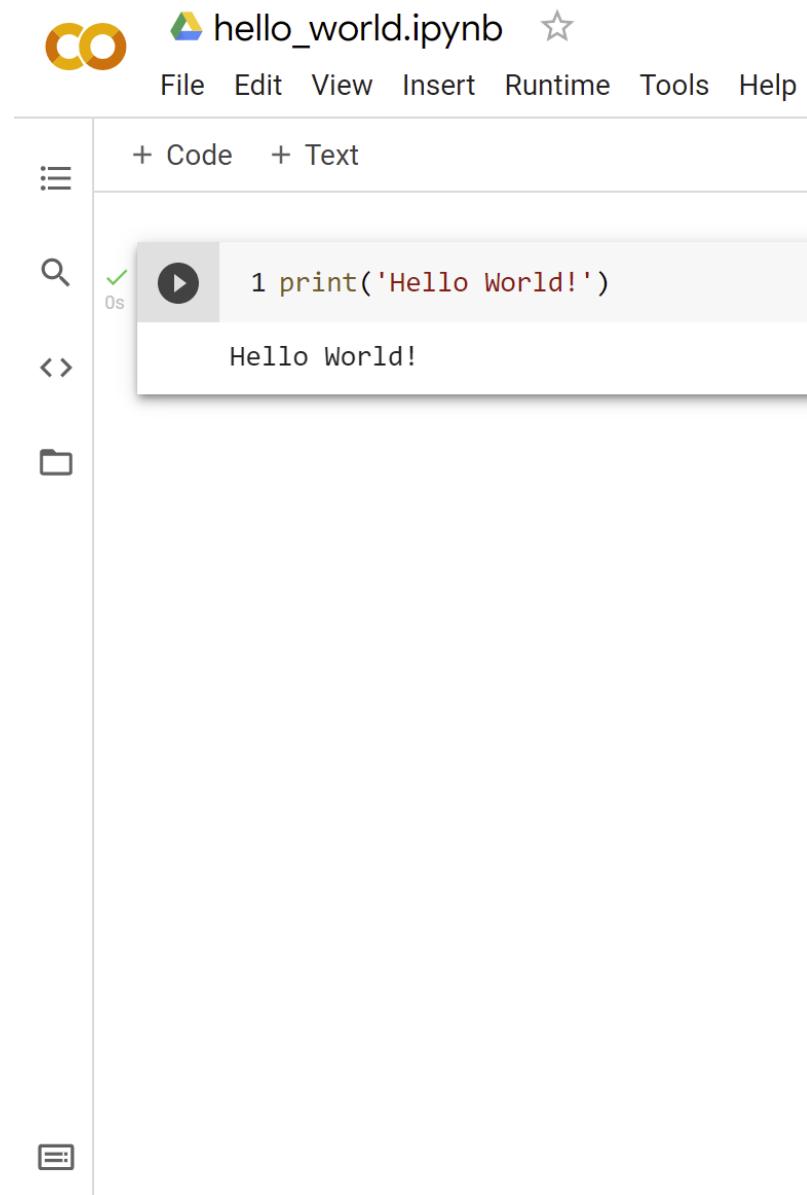
- 1-NN: for a query sample $x^{(q)}$, we predict the $y^{(q)}$ of the closest element in our dataset \mathcal{D}

$$y^{(q)} = y^{(p)}, \quad p = \arg \min_i (d(x^{(q)}, x^{(i)}))$$

- $d(\cdot)$ distance metric (usually Euclidean, can be learnt)
- Supervised, classification/regression, non-parametric (“lazy”, keeps the data instead of summarizing it into a model)
- In k-NN (usually k even), the k -nearest neighbours vote for the class
- Very simple and intuitive.
- It works reasonably well if the feature space is reasonably well sampled by \mathcal{D}
- Computationally demanding at test time!

Let's start coding!

- If you are not familiar with Python, **do the tutorial**
<https://docs.python.org/3/tutorial/>
- To make things easier at the beginning, use Google Colab
<https://colab.research.google.com/>
- Use sklearn documentation
<https://scikit-learn.org/> and google for help
- The ML community shares a lot of resources, enjoy!



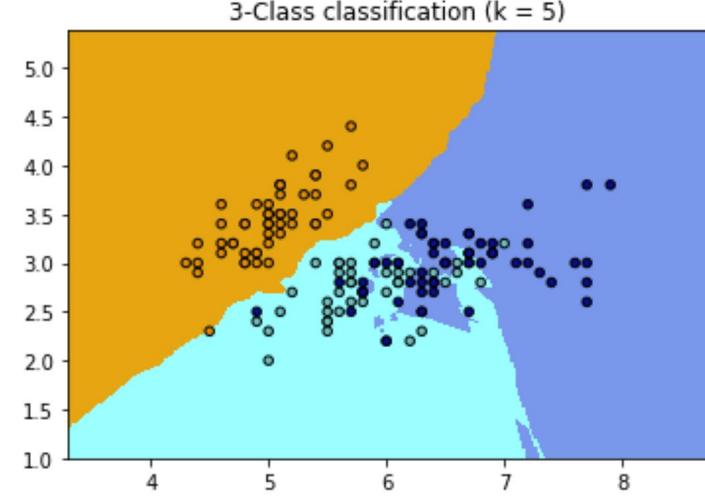
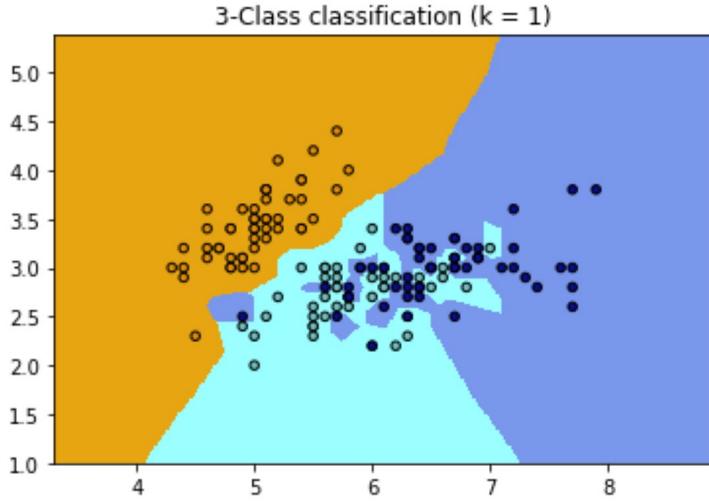
The screenshot shows the Google Colab interface. At the top, there's a file named "hello_world.ipynb" with a star icon. Below the header is a menu bar with File, Edit, View, Insert, Runtime, Tools, and Help. On the left, there's a sidebar with icons for code, text, search, and file navigation. The main area shows a code cell with the following content:

```
1 print('Hello World!')
```

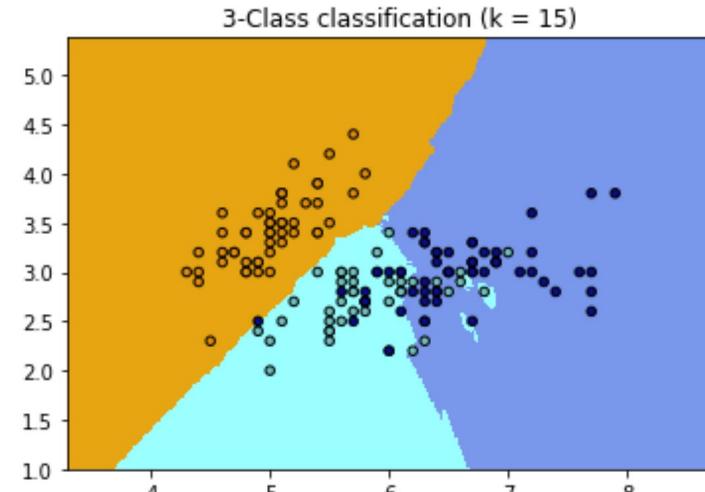
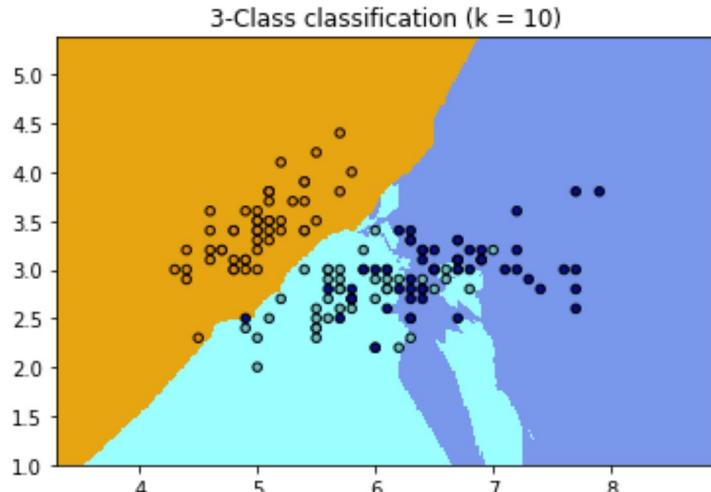
When run, the cell outputs "Hello World!".

K-nn examples – Iris

- Iris, sepal length and width, $k=\{1, 5, 10, 15\}$



Iris Setosa



Iris Versicolor

K-nn examples – MNIST

- Written digits

- Training size: 60K

k=1, accuracy=97.22%

k=3, accuracy=97.11%

k=5, accuracy=96.96%

k=10, accuracy=96.55%

k=15, accuracy=96.48%

k=20, accuracy=96.18%

k=25, accuracy=95.92%

k=30, accuracy=95.75%

- Training size: 5K

k=1, accuracy=93.61%

k=3, accuracy=93.12%

k=5, accuracy=93.08%

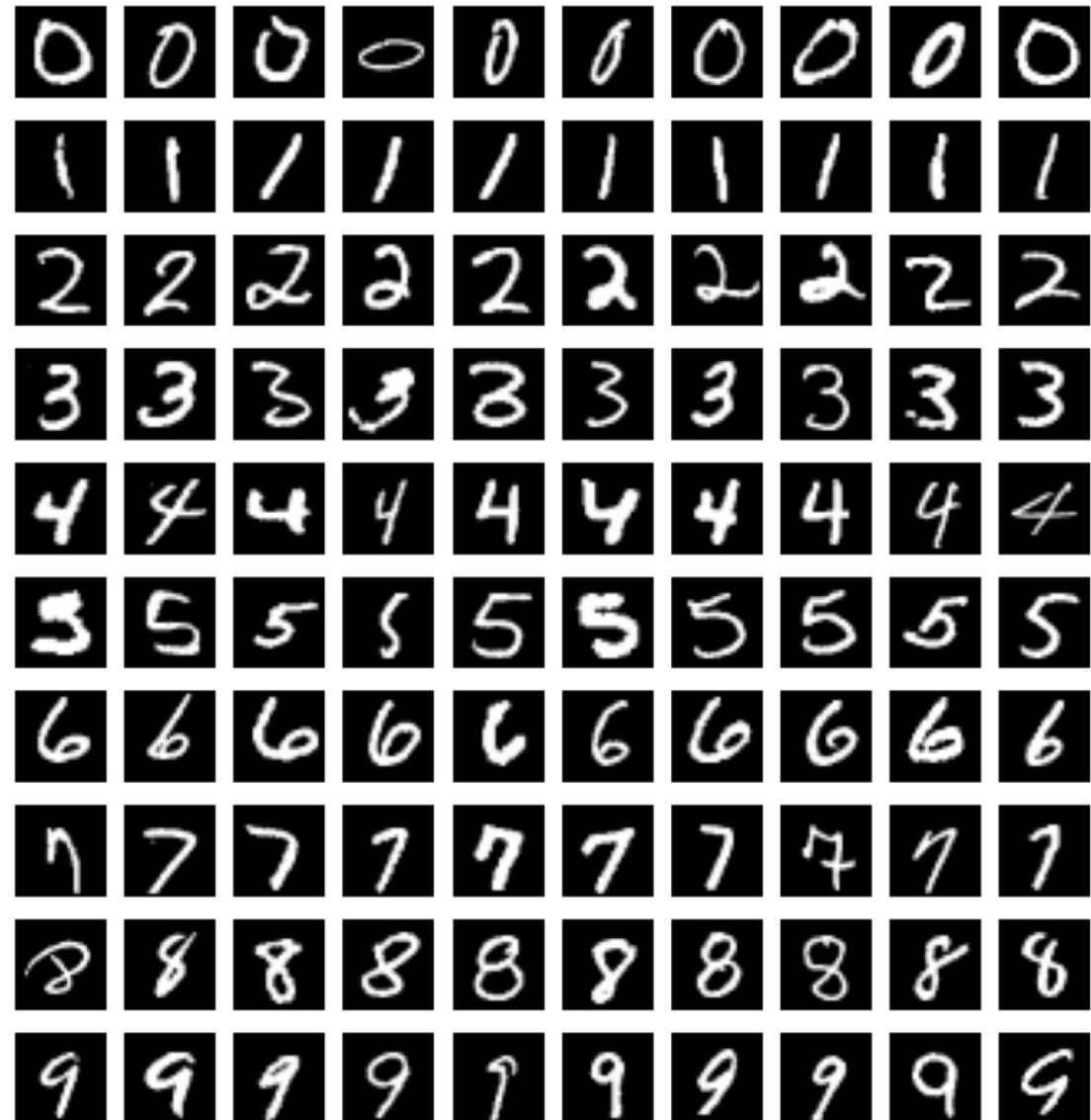
k=10, accuracy=92.37%

k=15, accuracy=91.89%

k=20, accuracy=91.45%

k=25, accuracy=90.93%

k=30, accuracy=90.46%



Colab link:

<https://colab.research.google.com/drive/1EdDg6rk23BW7rO9uwhogJf5YSMsZ-zK4>



K-nn examples – MNIST

- Changing the split:

 - Training size: 60K

k=1, accuracy=97.16%

k=3, accuracy=97.31%

k=5, accuracy=97.22%

k=10, accuracy=96.69%

k=15, accuracy=96.60%

k=20, accuracy=96.36%

k=25, accuracy=96.04%

k=30, accuracy=95.78%

 - Training size: 5K

k=1, accuracy=93.10%

k=3, accuracy=92.98%

k=5, accuracy=92.68%

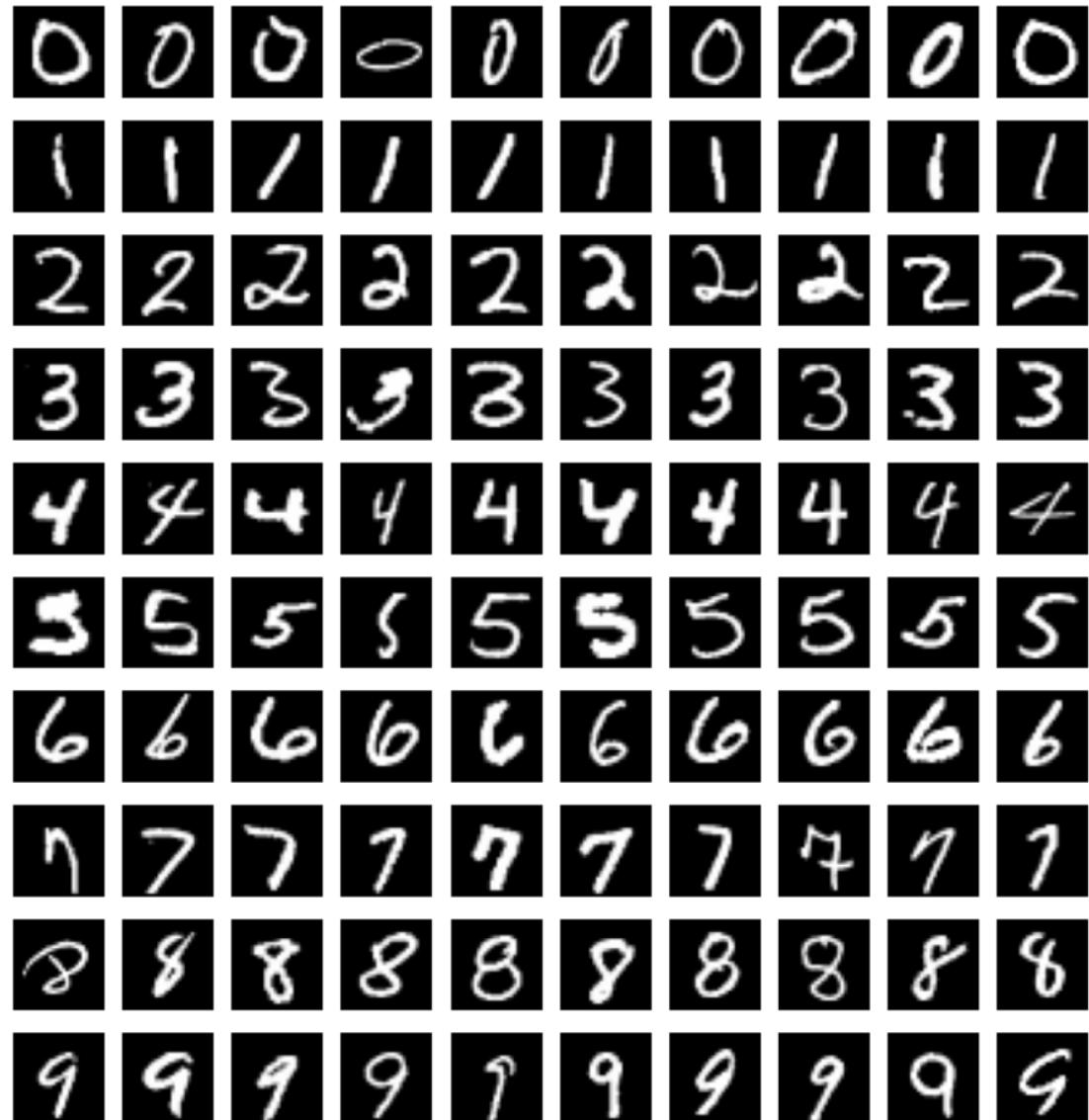
k=10, accuracy=92.14%

k=15, accuracy=91.60%

k=20, accuracy=90.90%

k=25, accuracy=90.37%

k=30, accuracy=89.97%



Colab link:

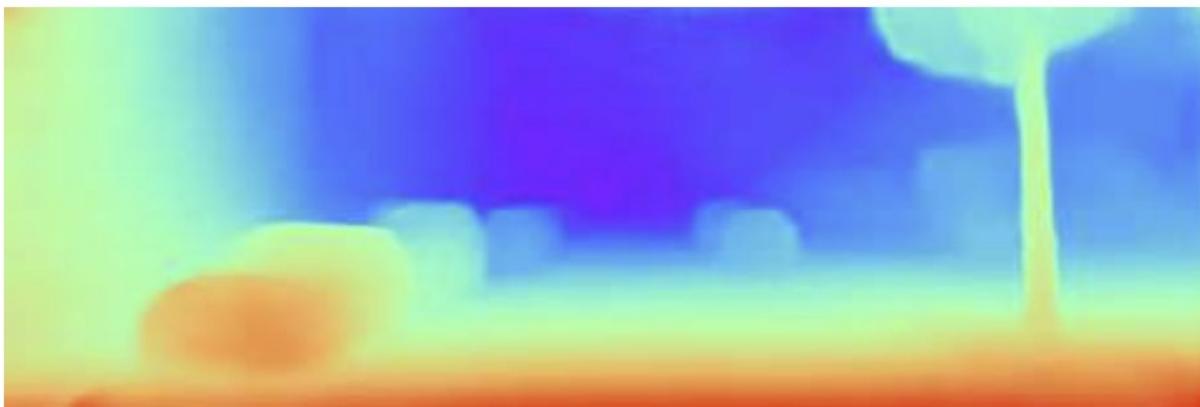
<https://colab.research.google.com/drive/1EdDg6rk23BW7rO9uwhogJf5YSMsZ-zK4>



Universidad
Zaragoza

Regression

- Learns a function that relates an input vector with a **continuous output**.
- Examples: depth prediction, weather forecasting, time-to-destination estimation...



Linear Regression

For more details, check Murphy chapter 11,
Deisenroth chapter 9

- The input/output relation is linear in the model parameters w

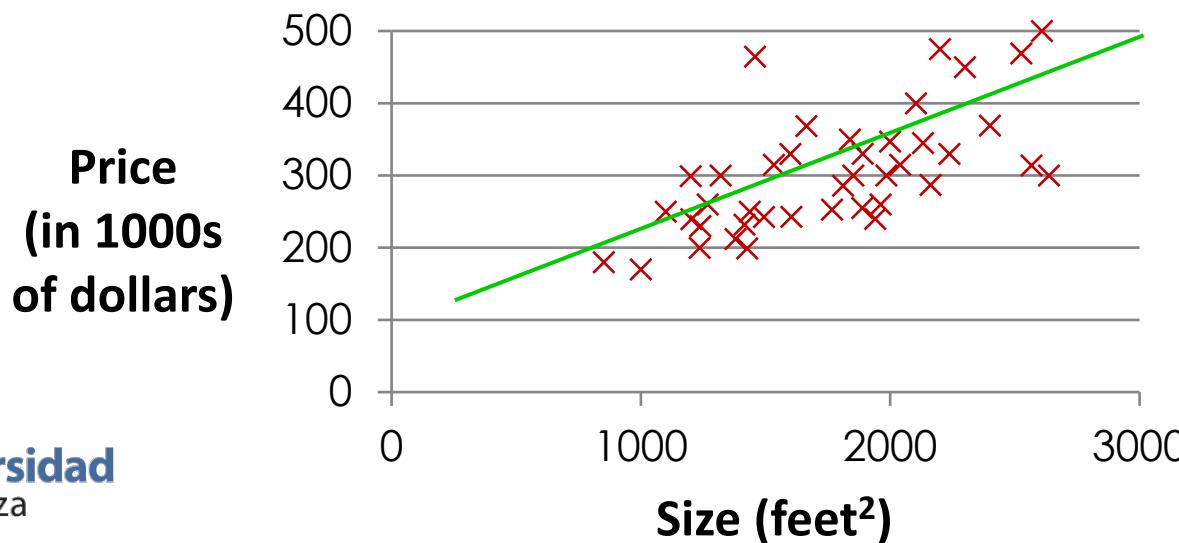
$$y = w^T x + \epsilon = \tilde{w}^T \tilde{x} + w_0 + \epsilon = \sum_{j=1}^d w_j x_j + \epsilon, \epsilon \sim N(0, \sigma^2)$$

$x = (1, x_1, \dots, x_d)^T \rightarrow$ independent/explanatory variables

$w = (w_0, w_1, \dots, w_d)^T, w_0 \rightarrow$ offset/bias, $w_{1:d} \rightarrow$ weights

$\tilde{x} = (x_1, \dots, x_d)^T, \tilde{w} = (w_1, \dots, w_d)^T$

- Example: Housing prices and size



How do we get w from data?

- Training data $\mathcal{D} = \{(x^{(1)}, y^{(1)}), \dots, (x^{(i)}, y^{(i)}), \dots, (x^{(N)}, y^{(N)})\}$

- Bayes' rule: $p(w|\mathcal{D}) = \frac{p(\mathcal{D}|w)p(w)}{p(\mathcal{D})}$

- MAP estimation:

$$\hat{w} = \arg \max_w p(w|\mathcal{D}) = \arg \max_w \frac{p(\mathcal{D}|w)p(w)}{p(\mathcal{D})} = \arg \max_w p(\mathcal{D}|w)p(w)$$

- MLE estimation (we do not have prior knowledge over w):

$$\hat{w} = \arg \max_w p(\mathcal{D}|w) = \arg \max_w p(y|x, w) = \arg \max_w \log(p(y|x, w))$$

- If the training samples are iid:

$$\hat{w} = \arg \max_w \log(p(y|x, w)) = \arg \max_w \log \left(\prod_i p(y^{(i)}|x^{(i)}, w) \right)$$

$$= \arg \max_w \sum_i \log(p(y^{(i)}|x^{(i)}, w))$$

How do we get w ?

- As $\epsilon \sim N(0, \sigma^2)$

$$\begin{aligned} & \arg \max_w \sum_i \log(p(y^{(i)} | x^{(i)}, w)) = \\ &= \arg \max_w \sum_i \log\left(\exp\left(-\frac{1}{2\sigma^2}(y^{(i)} - w^T x^{(i)})^2\right)\right) = \\ &= \arg \min_w \sum_i (y^{(i)} - w^T x^{(i)})^2 = \arg \min_w \sum_i r^{(i)^2} \end{aligned}$$

How do we get w ?

- We can define then a cost function

$$J(w) = \frac{1}{2} \sum_i (y^{(i)} - w^T x^{(i)})^2 = \frac{1}{2} \sum_i r^{(i)2}$$

- Our estimation \hat{w} is the one minimizing the cost function

$$\hat{w} = \arg \min_w J(w)$$

- Analytical solution for linear regression (ordinary least squares)

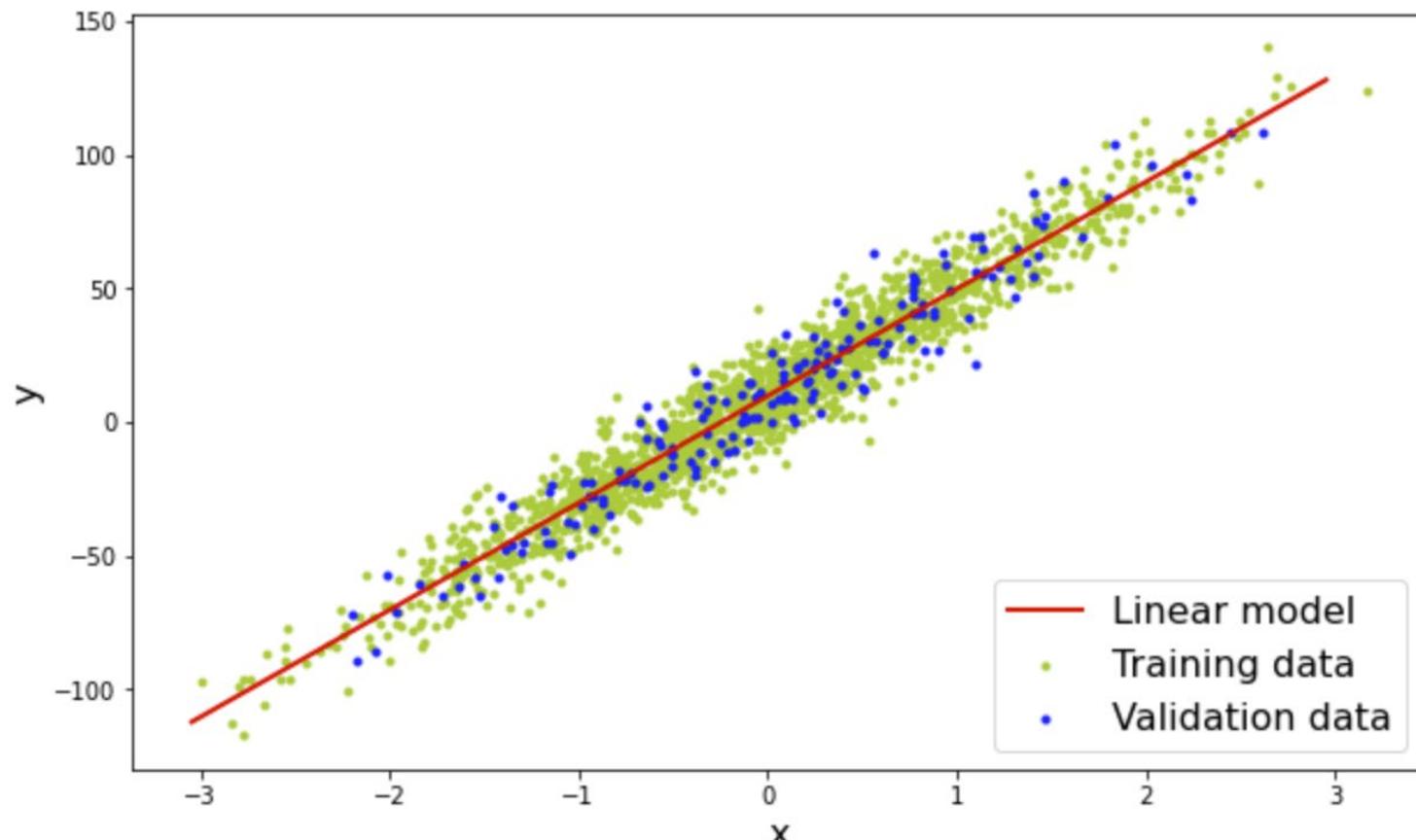
$$J(w) = \frac{1}{2} \sum_i (y^{(i)} - w^T x^{(i)})^2 = \frac{1}{2} (Xw - y)^T (Xw - y)$$
$$\rightarrow \dots \rightarrow \hat{w} = (X^T X)^{-1} X^T y$$

$$X = \begin{pmatrix} x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(N)} & \dots & x_d^{(N)} \end{pmatrix}, y = \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{pmatrix}$$

Simple linear regression example

- Colab link:

<https://colab.research.google.com/drive/1EXDJHTIvp4zK3qavGsXO03NMwYRr59JM>



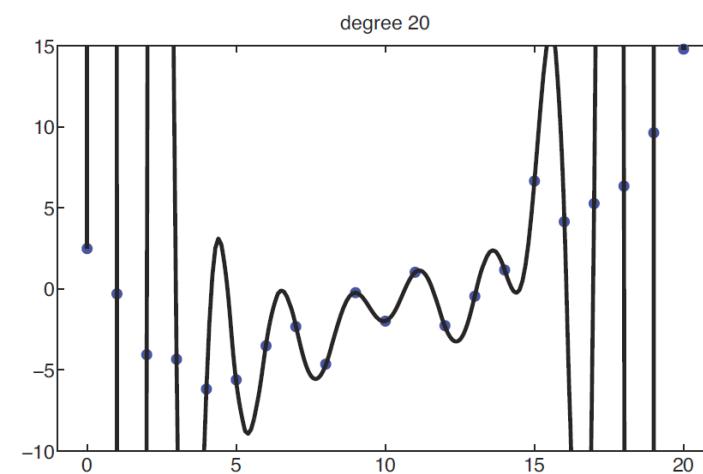
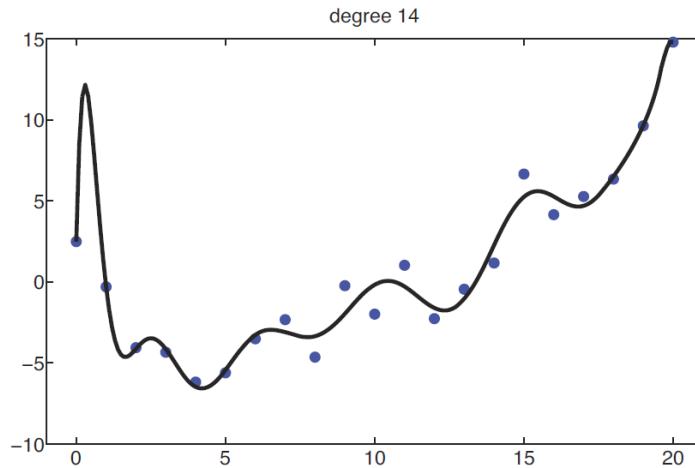
Nonlinear Regression (linear in w_i !!)

■ Basis function expansion

- Linear regression can be made to model non-linear relations by replacing $x \rightarrow \phi(x)$

■ Example:

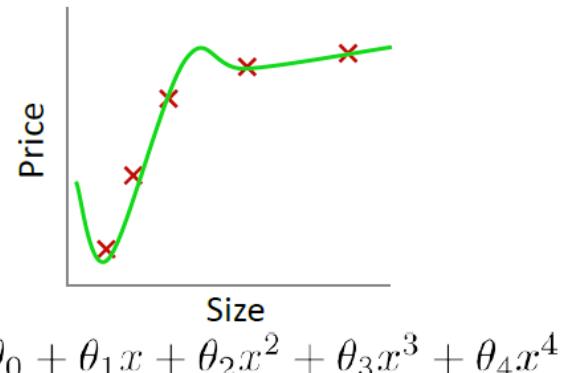
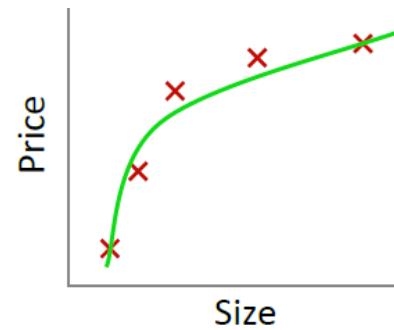
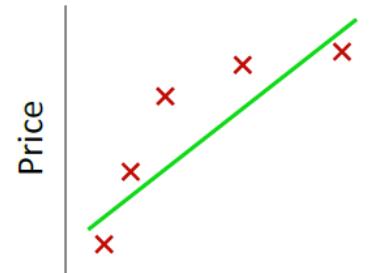
- Polynomial basis function $\phi(x) = [1, x, x^2, \dots, x^d]^T$ ($x \in \mathbb{R}^1, \phi(x) \in \mathbb{R}^d$)
- $y = w^T \phi(x) + \epsilon = \sum_{j=1}^d w_j x^j + \epsilon = w_0 + w_1 x^1 + w_2 x^2 + \dots + w_d x^d + \epsilon$
- Polynomials of degree 14 and 20 fit by LS (which one is better?)



Example and figures taken from Murphy's book

Overfitting

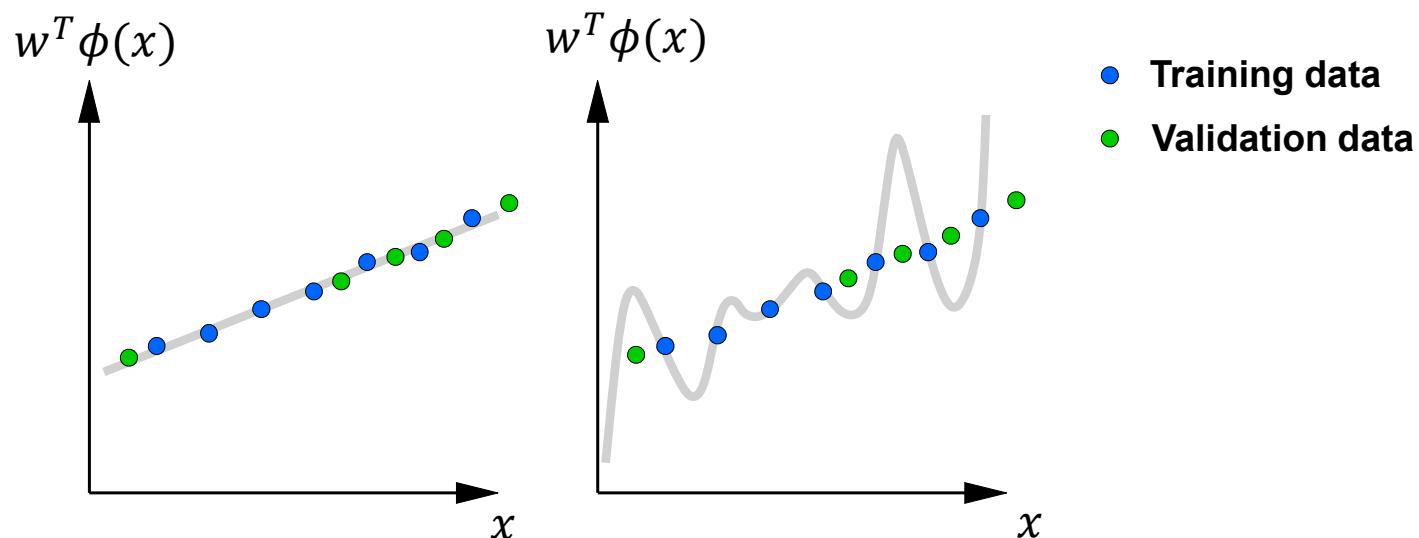
- The function $w^T \phi(x)$ fits the noise in addition to the underlying patterns. It “memorizes” the data instead of “learning to generalize”.
- The function does not capture the structure of the data.
- The predictive capability of an overfitted model is poor.
- Techniques:
 - Cross-validation
 - Regularization



$$J(\theta) = \frac{1}{2} \sum_{i=1}^N (h_\theta(x^{(i)}) - y^{(i)})^2 \approx 0$$

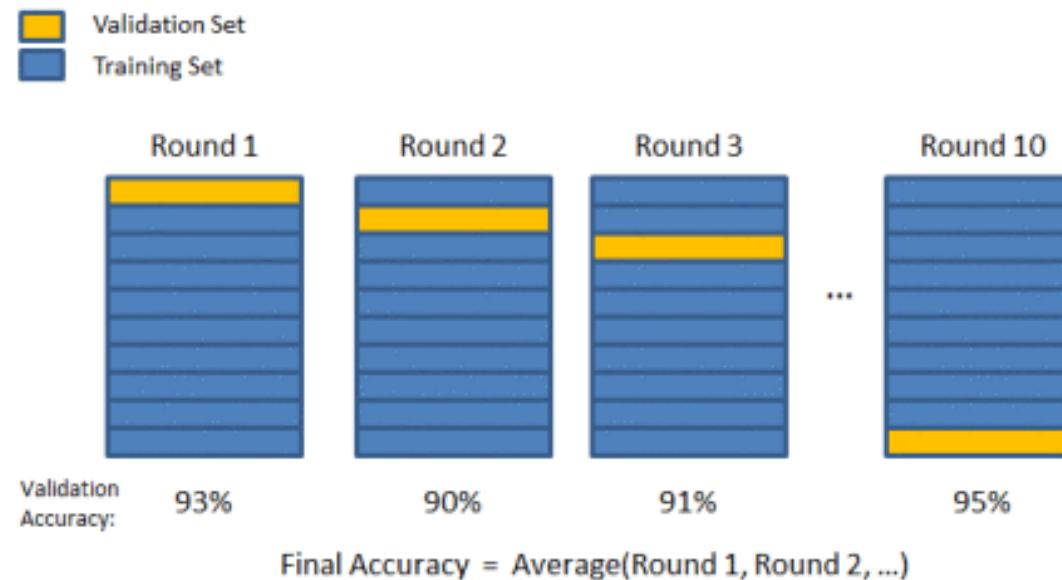
Cross-validation

- The data is partitioned into two sets.
 - The **training set** is used to estimate the model parameters w
 - The **validation set** is used to decide among a set of hyperparameters (models)
 - What we care about is generalization! → Final performance metrics should be given for a **test set**



K-fold cross-validation

- Common practice to have a higher degree of invariance to the validation/test split.
 - Separate the data into k folds
 - Do k training-validation rounds. A usual value for k is 10.
 - The validation/training sets are changed each round.
 - Metrics are averaged over all folds.



Summary: data splits

Model selection and fitting



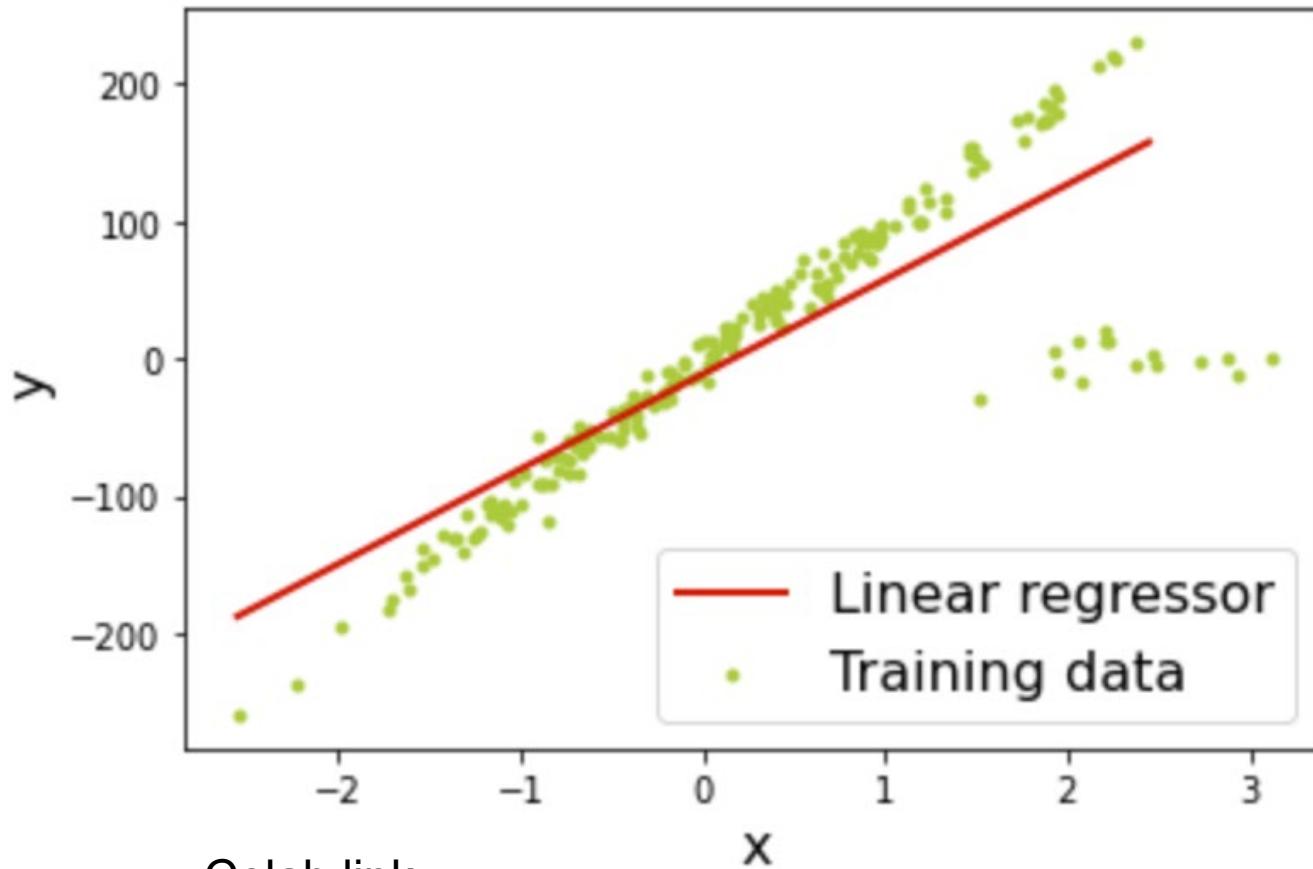
Test characterization



- Training data (~80% of the data)
 - Used to learn the model parameters
- Validation (~10% of the data)
 - Used to select the hyper-parameters
 - After used, can be added to the training data to train the final model
 - Training/validation splits can change during training (k-fold cross-validation)
- Test (~10% of the data)
 - Never used at the training stage!
 - Used to obtain the final specs of our model
 - After used, can be added to the training data to train the final model

Robust Regression

- The quadratic error penalty from the Gaussian noise gives a high weight to outliers

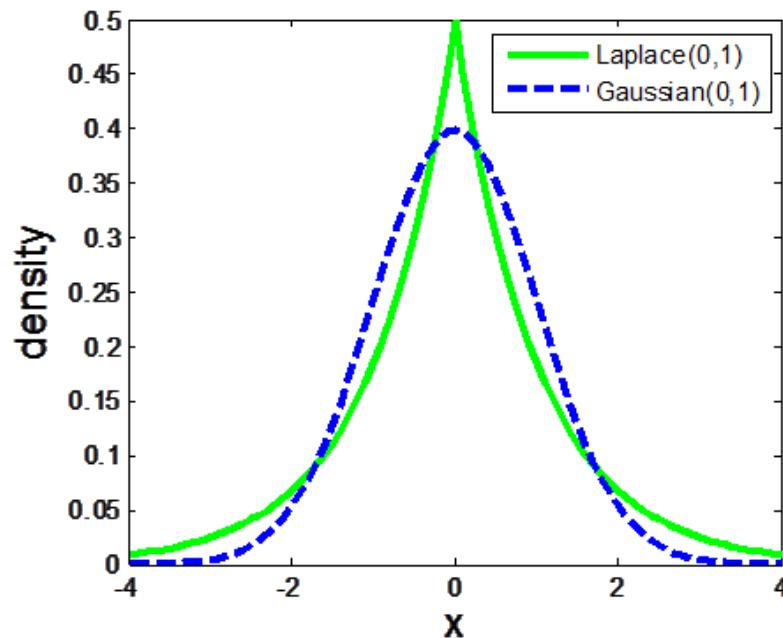


Colab link:

<https://drive.google.com/file/d/1G5E6QHihfK1TyaYe6dxlF2FF5XC1kXn->

Robust Regression

- Probabilistic view: The low tolerance of Gaussian pdfs for high-noise data samples causes problems with outliers
- One possible solution: assume longer-tailed distributions, e.g., the Laplace one
- Other: detect outliers and remove them
 - You will learn about RANSAC in the computer vision course



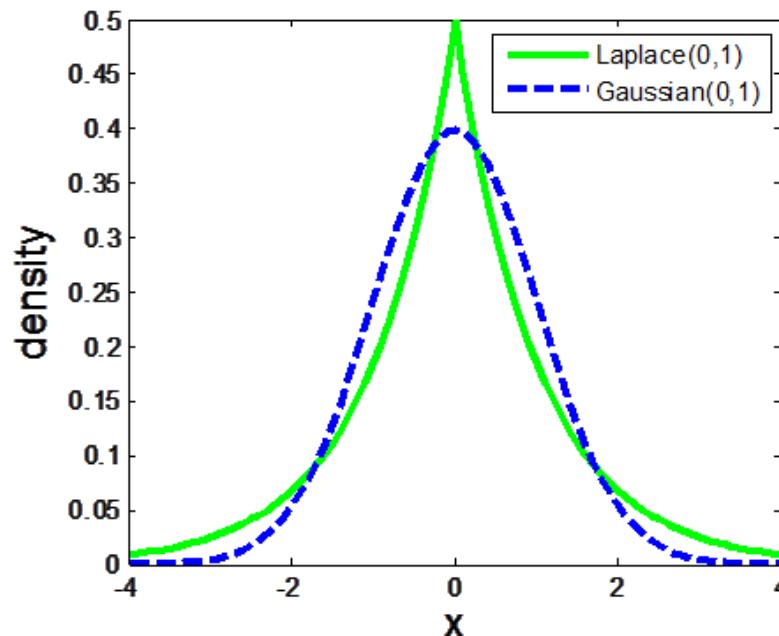
Robust Regression

- Laplace distribution:

$$p(y^{(i)}|x^{(i)}, w) = \exp\left(-\frac{1}{2b^2}|y^{(i)} - w^T x^{(i)}|\right) \rightarrow \dots \rightarrow$$

$$J(w) = \frac{1}{2} \sum_i |y^{(i)} - w^T x^{(i)}| = \sum_i |r^{(i)}|$$

- Non-linear, hard to optimize!

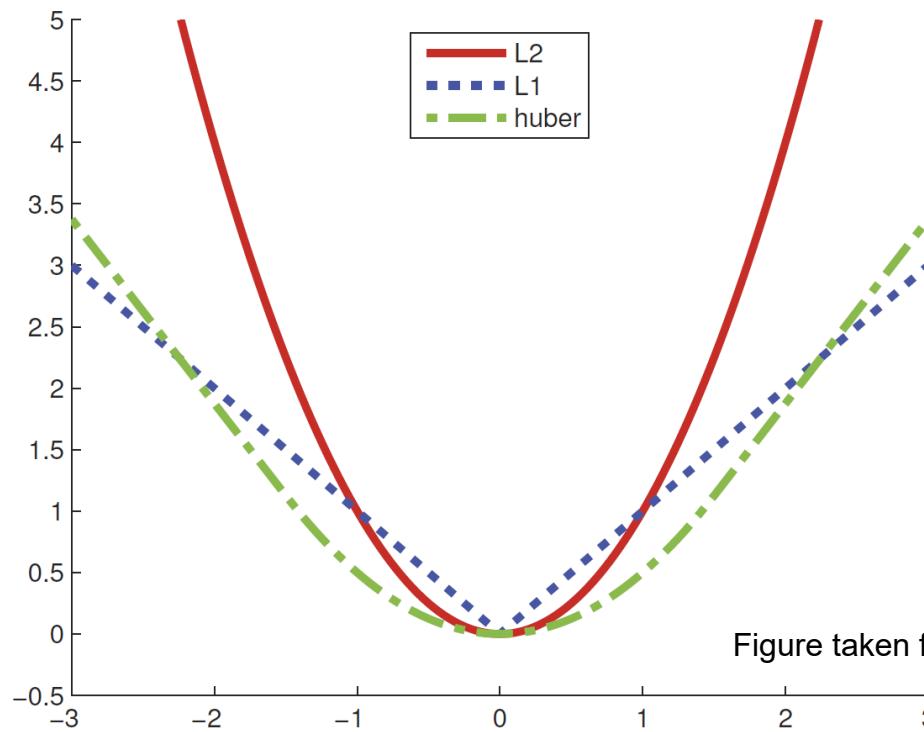


Robust Regression: Huber loss function

- Equivalent to L2 near zero, to L1 for large errors

$$r^{(i)} = \begin{cases} \frac{r^{(i)2}}{2} & \text{if } |r| \leq \delta \\ \delta|r| - \frac{\delta^2}{2} & \text{if } |r| > \delta \end{cases}$$

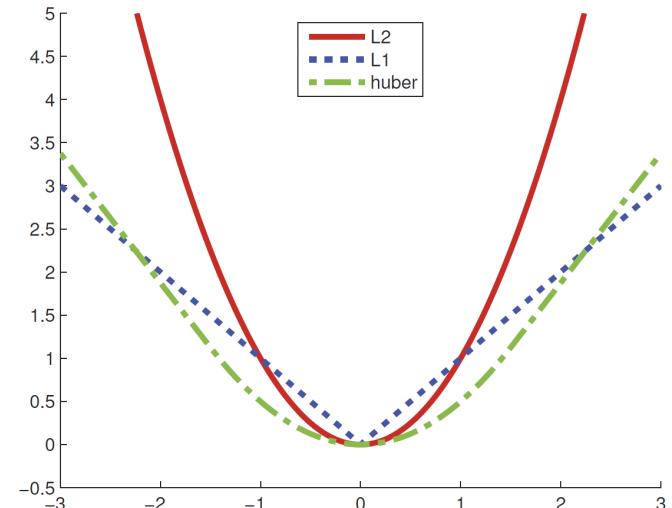
- Differentiable!



Huber loss function

- How do we tune the parameter δ ?

$$r^{(i)} = \begin{cases} \frac{r^{(i)2}}{2} & \text{if } |r| \leq \delta \\ \delta|r| - \frac{\delta^2}{2} & \text{if } |r| > \delta \end{cases}$$

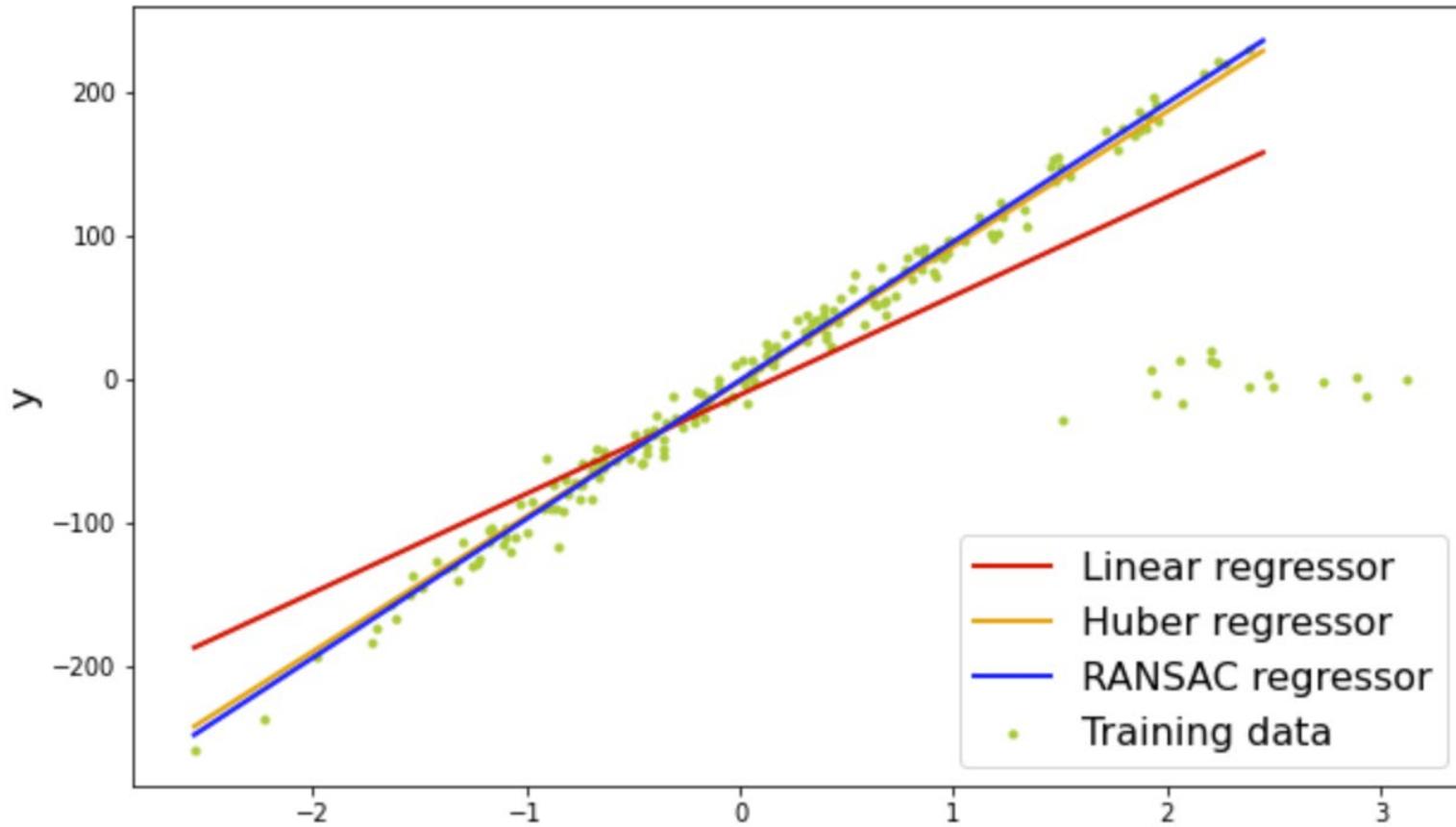


- δ is the “boundary” between inliers (Gaussian pdf assumed) and the outliers (Laplacian pdf assumed). So:

1. [Huber 1981], says: “good choices are between σ and 2σ ”
 - σ is the standard deviation of the residuals $r^{(i)}$ of the inliers
 2. If you do not know σ , you estimate it from the training data residuals $r^{(i)}$.
 - IMPORTANT! You have outliers, so use a robust estimator for σ . MAD (Median Absolute Deviation) is common:
- $$\hat{\sigma} = 1.4826 \cdot \text{MAD}, \text{MAD} = \text{median}(|r^{(1)} - \tilde{r}|, \dots, |r^{(N)} - \tilde{r}|), \tilde{r} = \text{median}(r^{(1)}, \dots, r^{(N)})$$

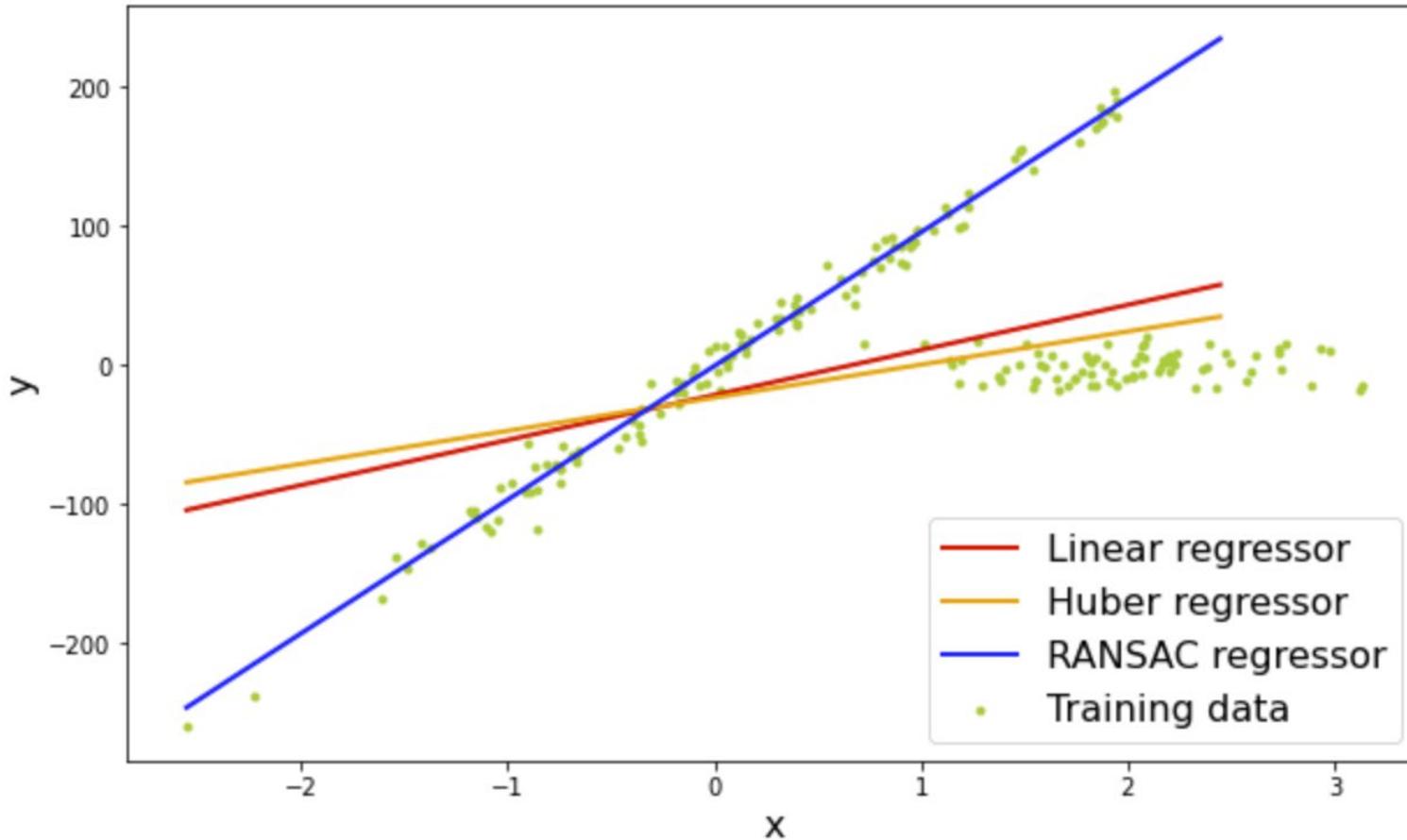
Robust regression example

- The Huber loss and RANSAC are robust to outliers



Robust regression example

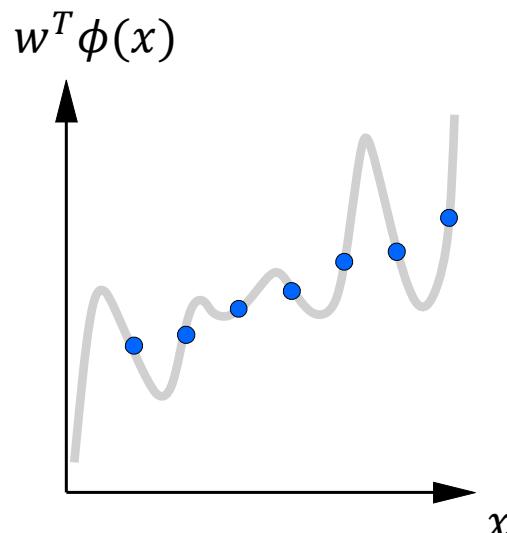
- But RANSAC can tolerate a higher outlier ratio!



Ridge Regression

- MLE can overfit easily for noisy training data
- Complex models will fit the noise
- Solution: Use a Gaussian prior over the model parameters (remember MAP from slide 23!)

$$p(w) \sim N(0, \tau^2)$$



Ridge Regression

- The MAP estimation is

$$\begin{aligned}\hat{w} &= \arg \max_w p(\mathcal{D}|w)p(w) = \\ &= \arg \max_w \sum_i \log \left(\exp \left(-\frac{1}{2\sigma^2} (y^{(i)} - w^T x^{(i)})^2 \right) \right) + \sum_j \log \left(\exp \left(-\frac{1}{2\tau^2} (w_j)^2 \right) \right)\end{aligned}$$

- Resulting in the following cost function

$$J(w) = \frac{1}{N} \sum_i (y^{(i)} - w^T x^{(i)})^2 + \lambda w^T w; \quad \lambda = \frac{\sigma^2}{\tau^2}$$

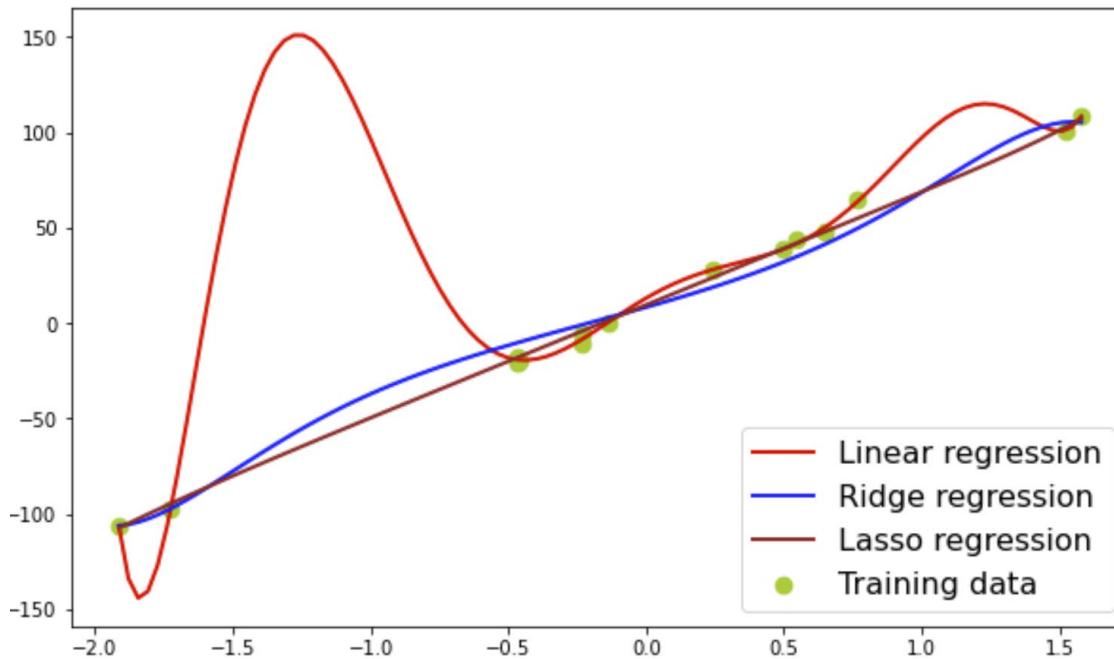
- Which has the following solution $\hat{w} = (\lambda I_d + X^T X)^{-1} X^T y$
- λ usually selected by cross-validation

Ridge Regression

Colab link:

<https://drive.google.com/file/d/1EnrFqx8Ks0L3-Uy99IB6DJ4DdEHyP5pl>

- Avoids overfitting
- $\lambda w^T w = \lambda \|w\|_2$ AKA “weight decay” in deep learning
- $\lambda \|w\|_1$ LASSO (Least Absolute Shrinkage and Selection Operator)
 - Assumes a Laplacian distribution for the weights (then $\lambda = \frac{2\sigma^2}{b^2}$)
 - Can assign zero values to some coefficients (feature selection)

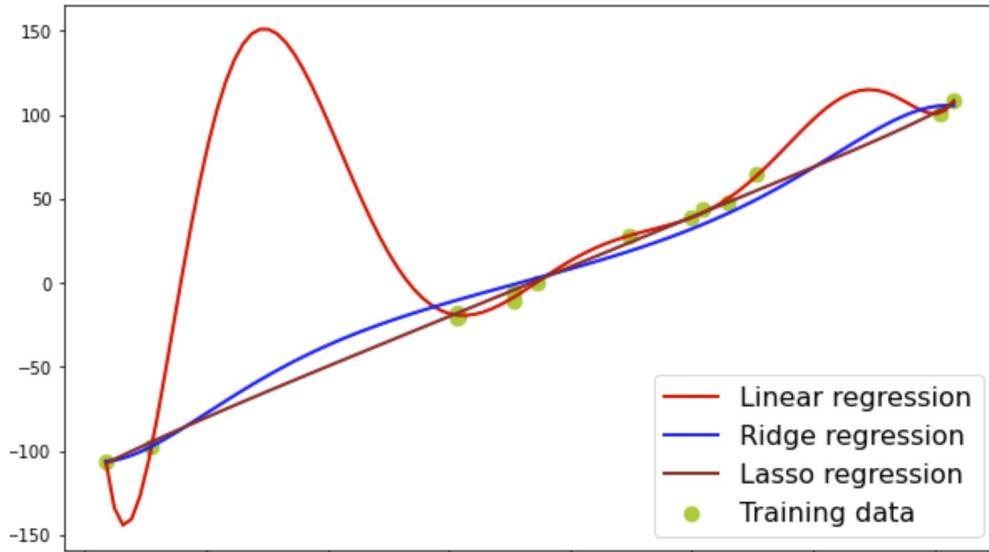


Data regularization

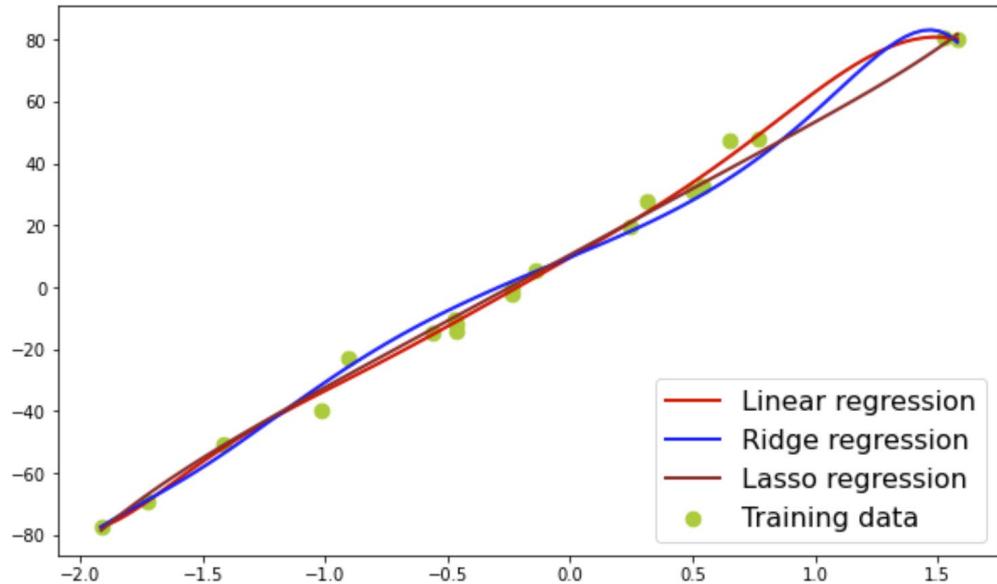
Colab link:

https://colab.research.google.com/drive/1j3kYQAHHeGSbEFf66U_f2P9wlSatT1Erl

- BTW, in this example, check the regularization effect of adding more data
 - 15 data samples

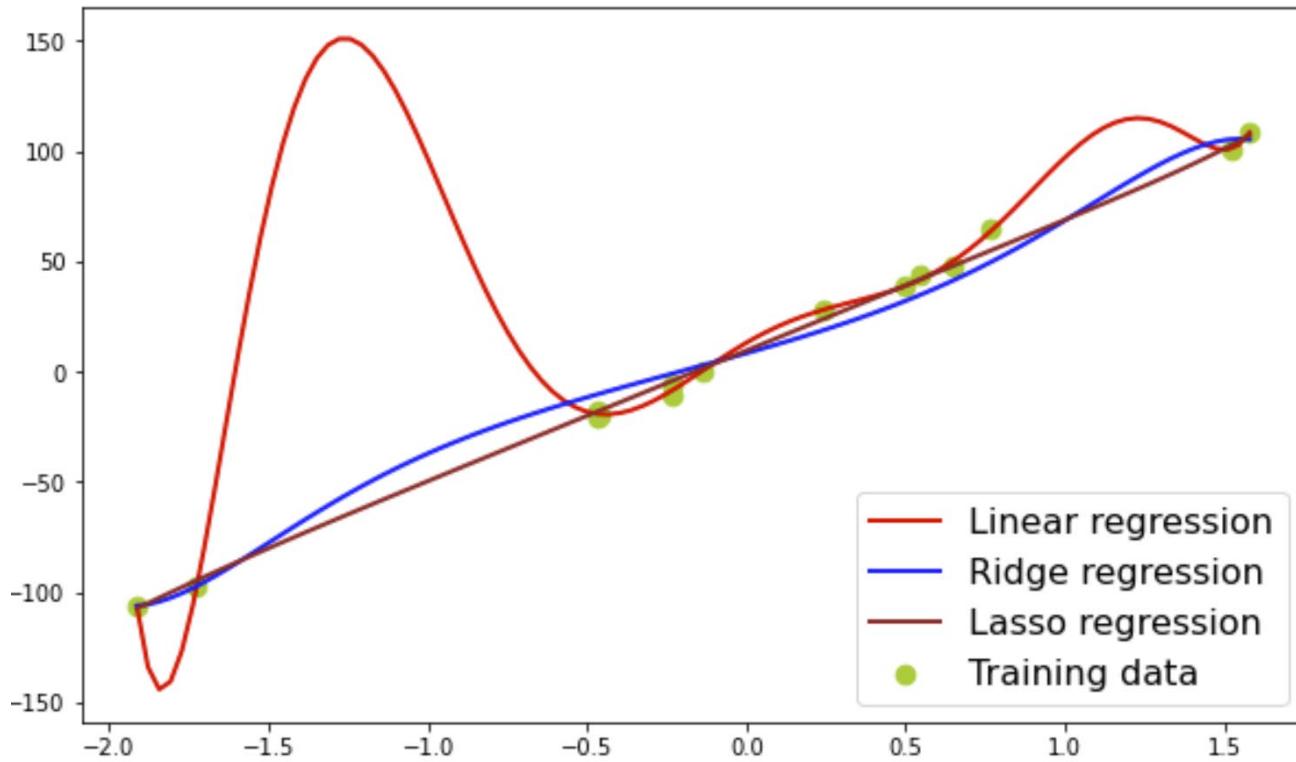


- 20 data samples



Summary: avoiding overfitting

- Do cross-validation to select the best model
- Add constraints on your model weights
- Add more training data!



Regression example: The diabetes dataset

■ Colab link:

<https://drive.google.com/file/d/1ElKbBCWetqEVKumrTX4GEpxdbRSFXwVw>

The screenshot shows the scikit-learn User Guide page for the Diabetes dataset. The top navigation bar includes links for 'Prev', 'Up', and 'Next'. Below it, a sidebar has links for 'scikit-learn 1.0' and 'Other versions', and a note to 'Please cite us if you use the software.' The main content area is titled '7.1.3. Diabetes dataset'. It describes the dataset as having ten baseline variables and one target variable, with 442 instances. The 'Data Set Characteristics' section provides details about the number of instances, attributes, target, and attribute information, listing all 11 columns of the dataset.

scikit
learn

Prev Up Next

scikit-learn 1.0 Other versions

Please cite us if you use the software.

User Guide

1. Supervised learning
2. Unsupervised learning
3. Model selection and evaluation
4. Inspection
5. Visualizations
6. Dataset transformations
7. Dataset loading utilities
- 7.1. Toy datasets
- 7.2. Real world datasets
- 7.3. Generated datasets
- 7.4. Loading other datasets
8. Computing with scikit-learn
9. Model persistence
10. Common pitfalls and recommended practices

Toggle Menu

7.1.3. Diabetes dataset

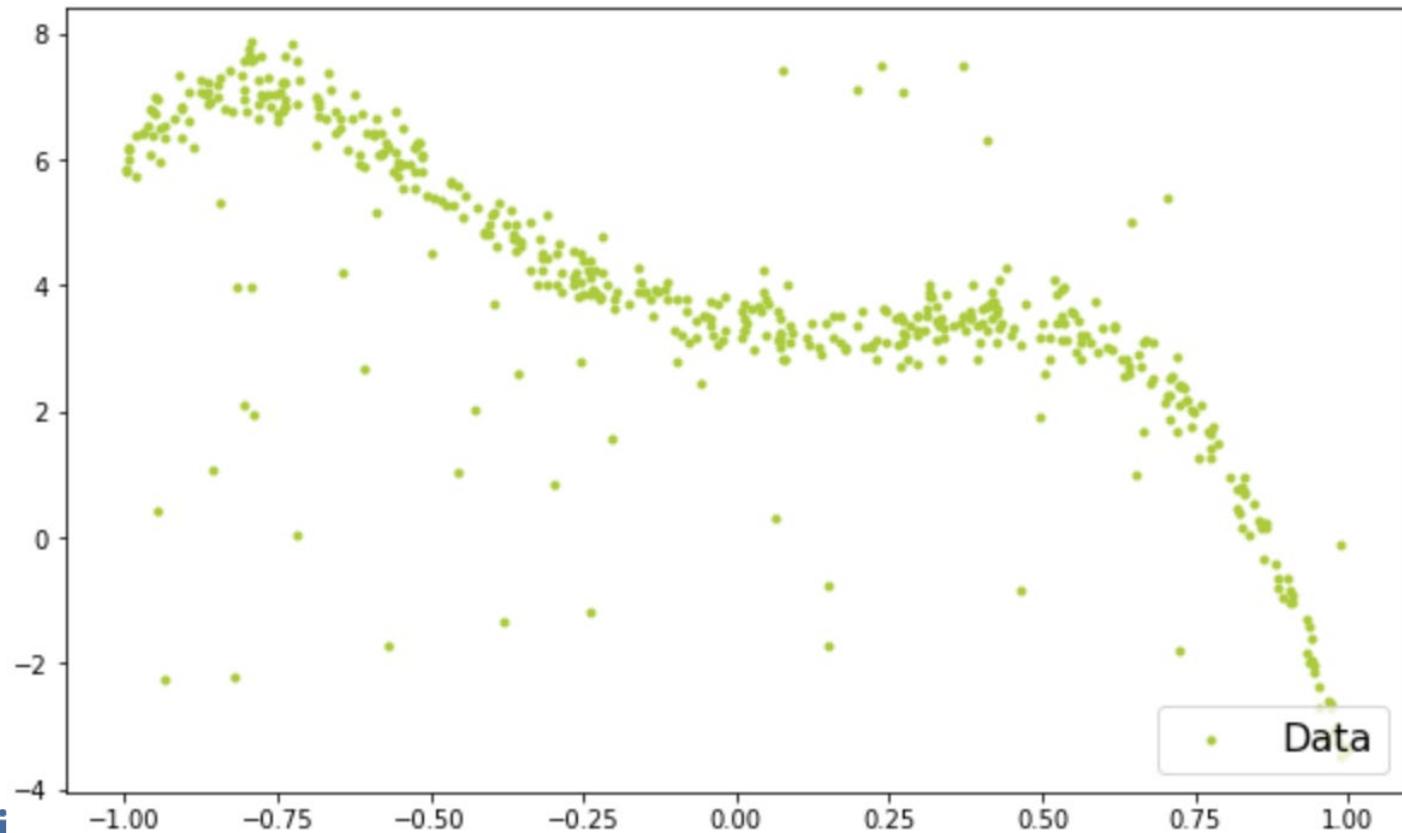
Ten baseline variables, age, sex, body mass index, average blood pressure, and six blood serum measurements were obtained for each of $n = 442$ diabetes patients, as well as the response of interest, a quantitative measure of disease progression one year after baseline.

Data Set Characteristics:

Number of Instances:	442
Number of Attributes:	First 10 columns are numeric predictive values
Target:	Column 11 is a quantitative measure of disease progression one year after baseline
Attribute Information:	<ul style="list-style-type: none">age age in yearssexbmi body mass indexbp average blood pressures1 tc, total serum cholesterols2 ldl, low-density lipoproteinss3 hdl, high-density lipoproteinss4 tch, total cholesterol / HDLs5 ltg, possibly log of serum triglycerides levels6 glu, blood sugar level

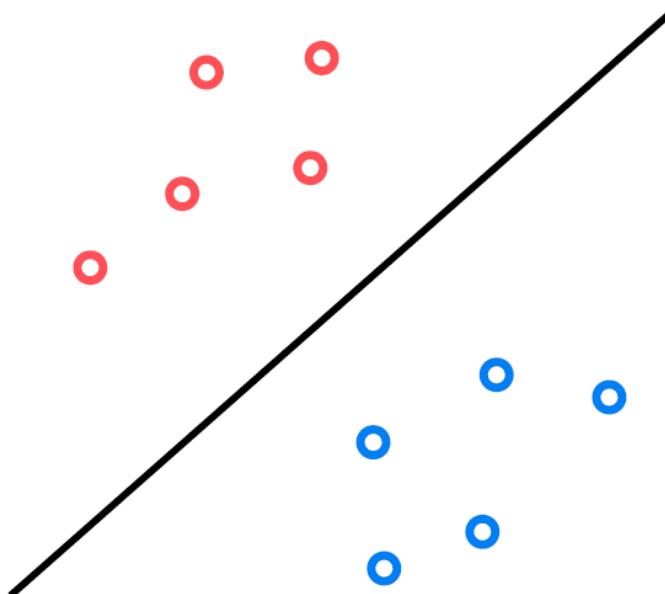
Exercise: Fit a polynomial model

- Link to Colab: <https://drive.google.com/file/d/1G1N9C9Yry-3Lc9A1mL5yawM1rDB6TnKM>
- Link to data: https://drive.google.com/file/d/1G_XseoB0lhII21jGdEkSYQyaEtMoSfQy



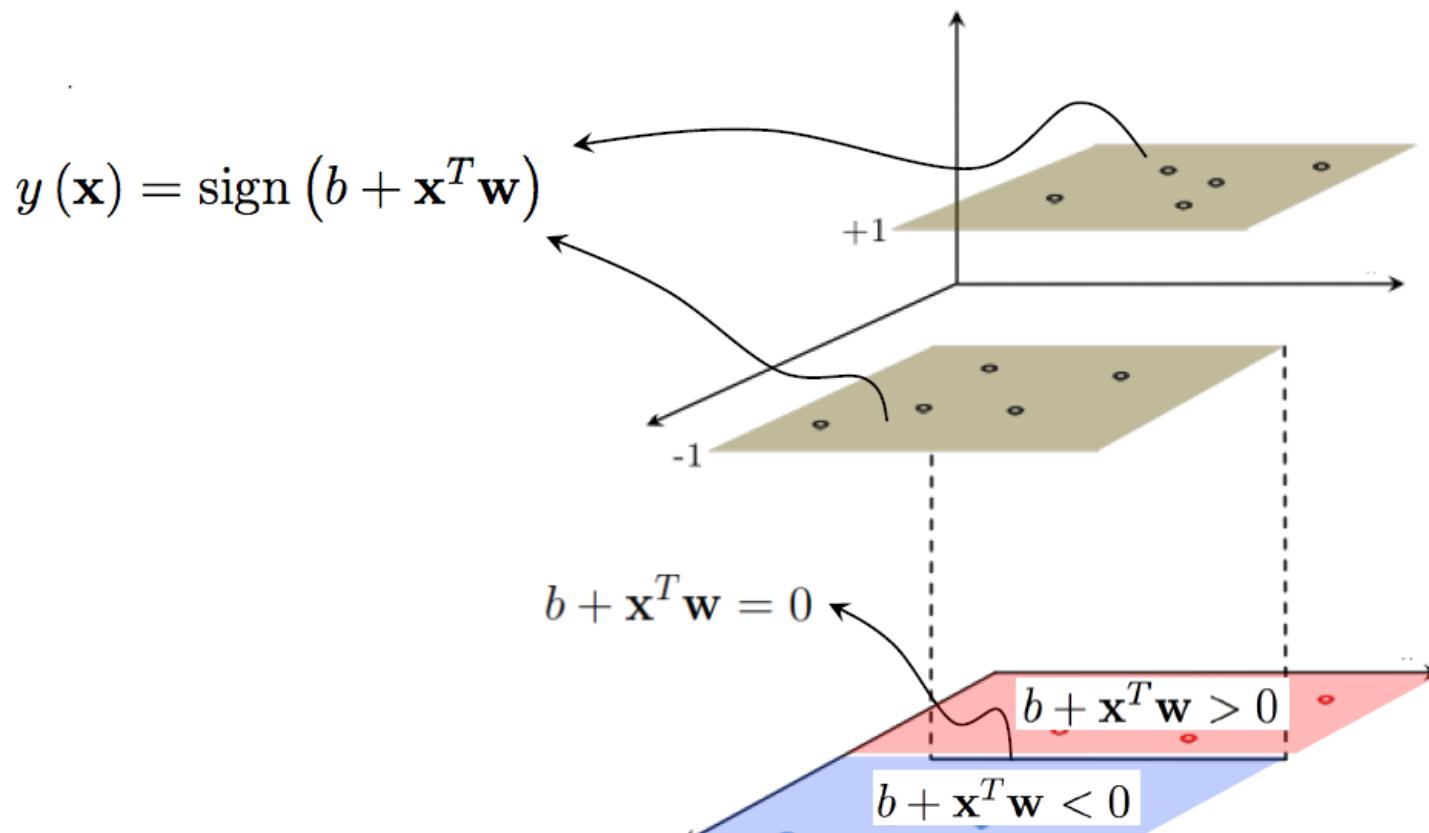
Linear Classification

- Classification: The output of the learned model is discrete
- Linear classification: The model is based on a linear combination of the input features
- Intuition: there is an hyperplane $\pi \in \mathbb{R}^{d-1}$ (e.g., a line for a 2D space) that separates the two classes/categories.



Binary Classification (intuition)

- We can get a discrete output $y \in \{-1, +1\}$ by taking the sign of the signed point-to-hyperplane distance
- The linear decision boundary is perpendicular to w
- The formulation can be extended easily to multi-class

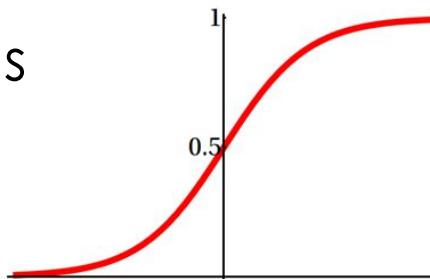


Sigmoid and tanh (intuition)

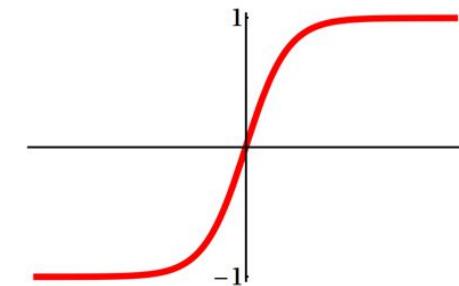
- The sign() function is not continuous.
- We can use continuous functions that approximate the sign function.
 - The sigmoid function
 - The tanh function

$$\sigma(\Sigma) = \frac{1}{1+e^{-\Sigma}}$$

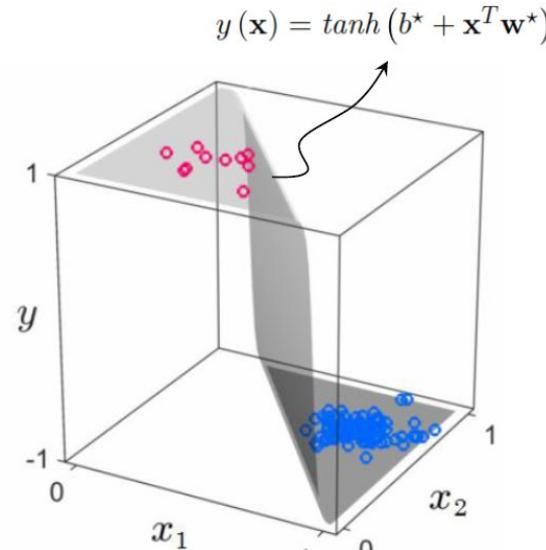
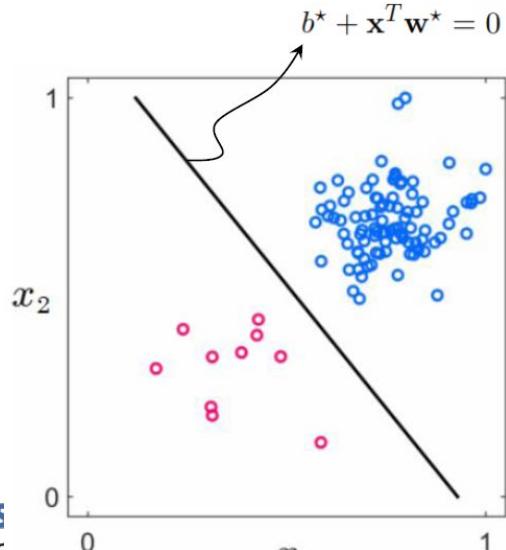
$$\tanh(\Sigma) = \frac{e^{\Sigma} - e^{-\Sigma}}{e^{\Sigma} + e^{-\Sigma}}$$



logistic (sigmoid, unipolar)



tanh (bipolar)



Logistic Regression (probabilistic interpretation)

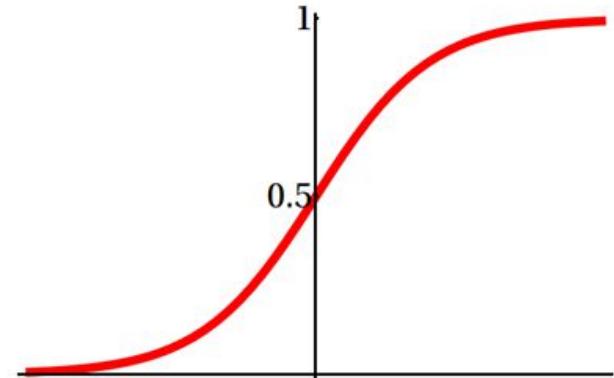
- The input/output model is

$$p(y|x, w) = \text{Ber}(y|\sigma(w^T x)) \rightarrow$$

$$\begin{cases} p(y = 1|x, w) = \sigma(w^T x) \\ p(y = 0|x, w) = 1 - \sigma(w^T x) \end{cases}$$

- $w^T x \rightarrow \infty \Rightarrow \sigma(w^T x) \rightarrow 1 \Rightarrow$
 $\Rightarrow p(y = 1|x, w) \rightarrow 1 \ (p(y = 0|x, w) \rightarrow 0)$
- $w^T x \rightarrow -\infty \Rightarrow \sigma(w^T x) \rightarrow 0 \Rightarrow$
 $\Rightarrow p(y = 1|x, w) \rightarrow 0 \ (p(y = 0|x, w) \rightarrow 1)$
- $w^T x = 0 \Rightarrow \sigma(w^T x) \rightarrow 0.5 \Rightarrow$
 $\Rightarrow p(y = 1|x, w) = 0.5 \ (p(y = 0|x, w) \rightarrow 0.5)$

$$\sigma(\Sigma) = \frac{1}{1 + e^{-\Sigma}}$$



logistic (sigmoid, unipolar)

For more details, check Murphy chapter 10

How do we get w ?

- Training data $\mathcal{D} = \{(x^{(1)}, y^{(1)}), \dots, (x^{(i)}, y^{(i)}), \dots, (x^{(N)}, y^{(N)})\}$
 - $x^{(i)} \in \mathbb{R}^d, y^{(i)} \in \begin{cases} \{0, 1\} \\ \{-1, +1\} \end{cases}$ (depends if we use sigmoid or tanh)
- Bayes' rule: $p(w|\mathcal{D}) = \frac{p(\mathcal{D}|w)p(w)}{p(\mathcal{D})}$
- MAP estimation (leads to L1/L2 norm regularization):
$$\hat{w} = \arg \max_w p(w|\mathcal{D}) = \arg \max_w \frac{p(\mathcal{D}|w)p(w)}{p(\mathcal{D})}\\ = \arg \max_w p(\mathcal{D}|w)p(w)$$
- MLE estimation (we do not have prior knowledge over w):
$$\hat{w} = \arg \max_w p(\mathcal{D}|w) = \arg \max_w p(y|x, w)\\ = \arg \max_w \log(p(y|x, w))$$

How do we get w ?

- If the training samples are iid:

$$\hat{w} = \arg \max_w \log(p(y|x, w)) = \arg \max_w \log\left(\prod_i p(y^{(i)}|x^{(i)}, w)\right)$$

$$= \arg \min_w - \sum_i \log(p(y^{(i)}|x^{(i)}, w))$$

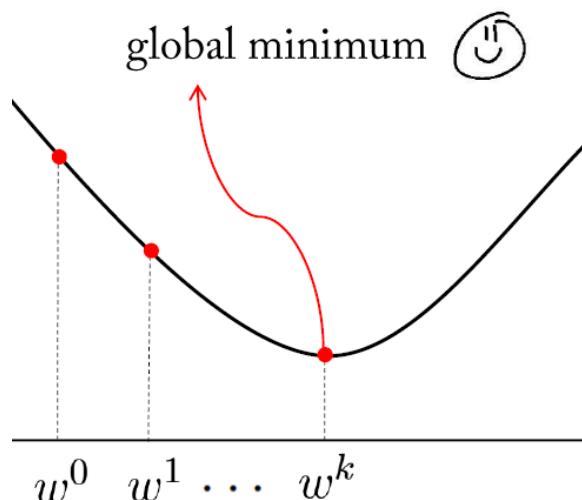
- $p(y^{(i)}|x^{(i)}, w)$ is the one described before, resulting in

$$\hat{w} = \arg \min_w - \sum_i \log(p(y^{(i)}|x^{(i)}, w)) = \dots =$$

$$= \begin{cases} - \sum_i (y^{(i)} \log(w^T x^{(i)}) + (1 - y^{(i)}) \log(1 - w^T x^{(i)})), & \text{if } y^{(i)} \in \{0, 1\} \\ - \sum_i (\log(1 + \exp(-y^{(i)} w^T x^{(i)}))), & \text{if } y^{(i)} \in \{-1, +1\} \end{cases}$$

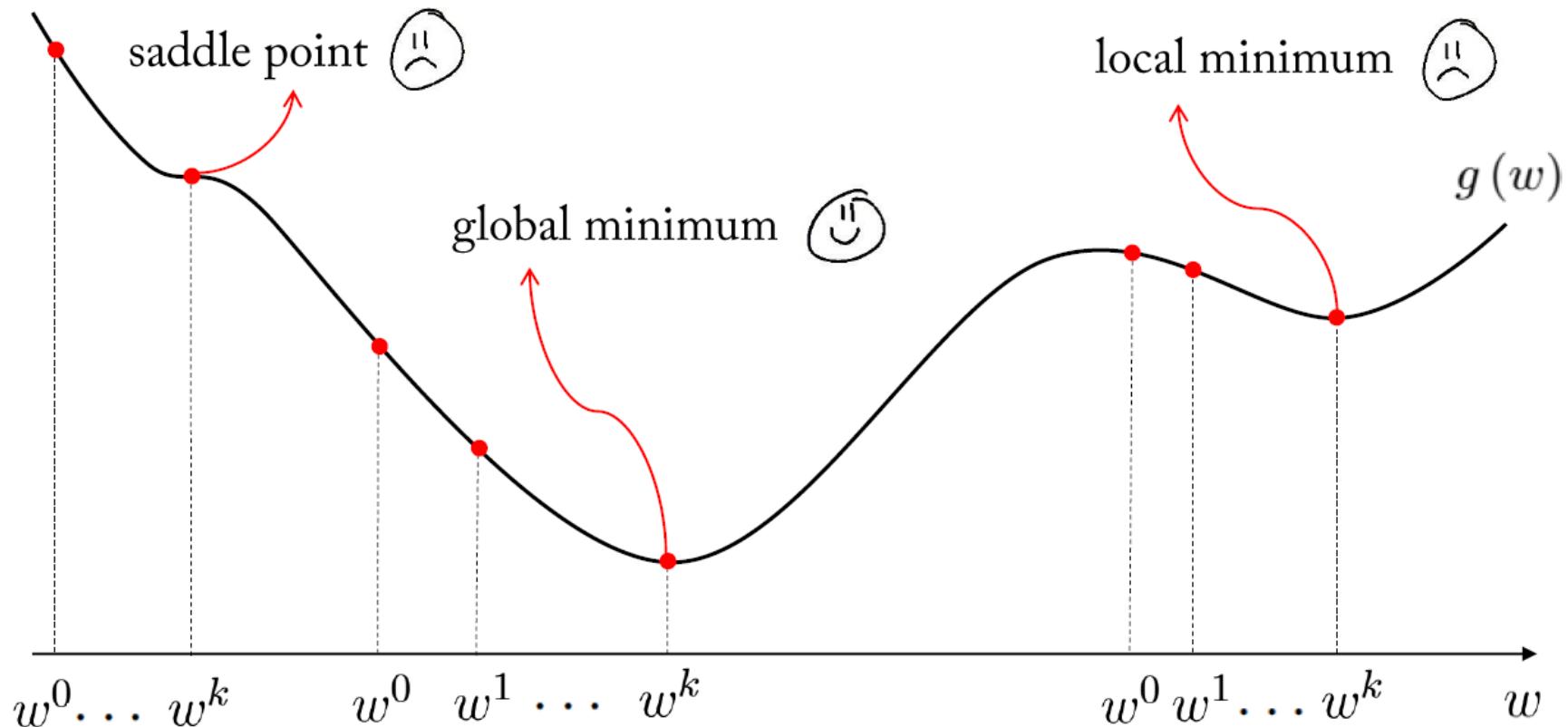
Cost function (loss)

- Cross-entropy loss: $L =$
$$\begin{cases} -\sum_i (y^{(i)} \log(w^T x^{(i)}) + (1 - y^{(i)}) \log(1 - w^T x^{(i)})), & \text{if } y^{(i)} \in \{0, 1\} \\ \sum_i (\log(1 + \exp(-y^{(i)} w^T x^{(i)}))), & \text{if } y^{(i)} \in \{-1, +1\} \end{cases}$$
- Differently from linear regression, no closed-form solution!...
- ... but, the loss is convex and has a single minimum (gradient descent will take us to it)



Gradient descent

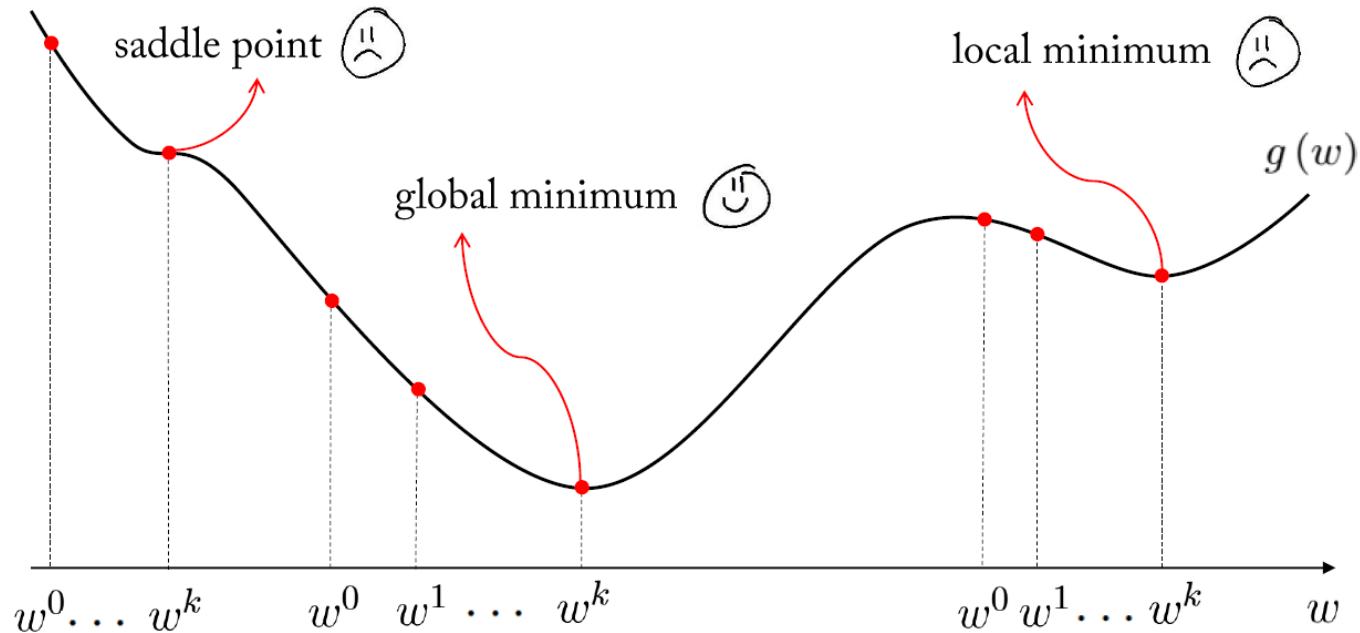
- Gradient descent → Iterative weight updates in the downward direction, $w_{k+1} = w_k - \eta_k \frac{dg(w)}{d(w)}$
- In general, it converges to a local minimum



Gradient descent

■ Issues:

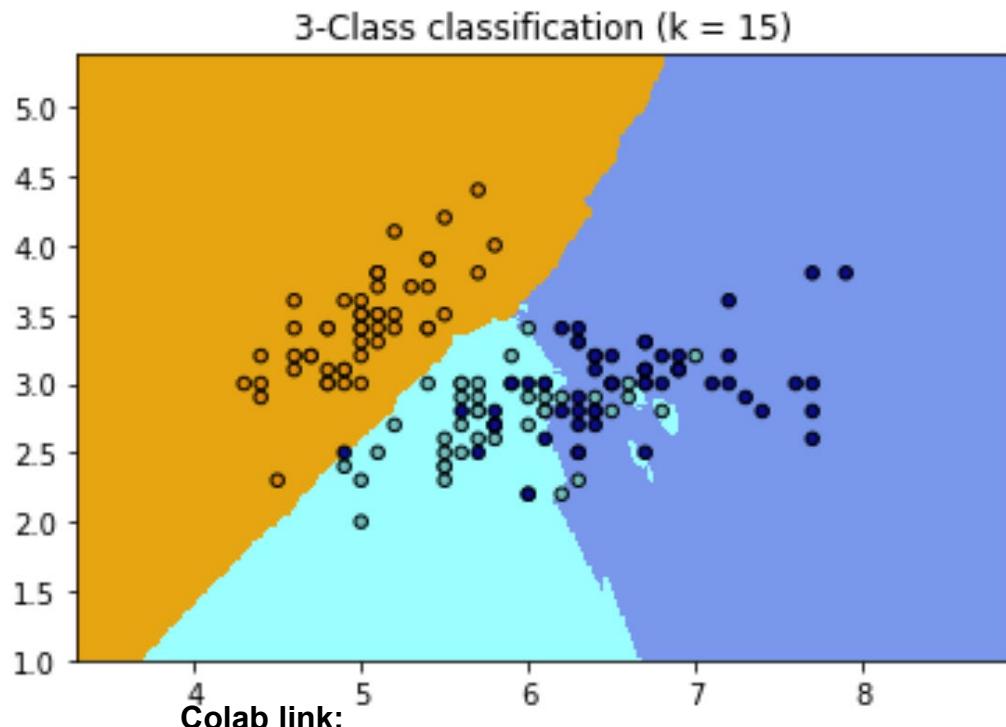
- Select the initial seed w_0
- Select the learning rate η_k
 - Small $\eta_k \rightarrow$ slow convergence, large $\eta_k \rightarrow$ fail to converge



Logistic regression in the Iris dataset

■ Take-home messages:

- Slightly outperforms k-nn (remember from our plots in slide 18, class boundaries are quite linear)
- More sophisticated methods (SVM) do not improve, again class boundaries are quite linear
- k-fold cross-validation IS relevant!



Iris Setosa



Iris Versicolor



Iris Virginica

Data standardization

- Many ML algorithms (SVM, L2/L1 regularization) assume all dimensions of x has zero-mean and similar covariance
- Most usual standardization:

$$x = (1, x_1, \dots, x_d)^T \rightarrow x' = (1, x'_1, \dots, x'_d)^T$$
$$x'_j = \frac{x_j - \mu_j}{\sigma_j}$$

- For more details, you can check <https://scikit-learn.org/stable/modules/preprocessing.html>
 - For example, for data with outliers, you should use robust estimates for mean and covariance

Ensembling

- The machine learning equivalent to the wisdom of the crowd
 - Aggregating predictions from several models improve over the individual ones
- Example (from Géron 2019):
 - Aggregating the votes of 1,000 binary classifiers with individual accuracies of 51% gives us an accuracy of 75%!
- **Practical trick:** when, at the end of a machine learning project, you trained several good models, combine their predictions to obtain an even better one.

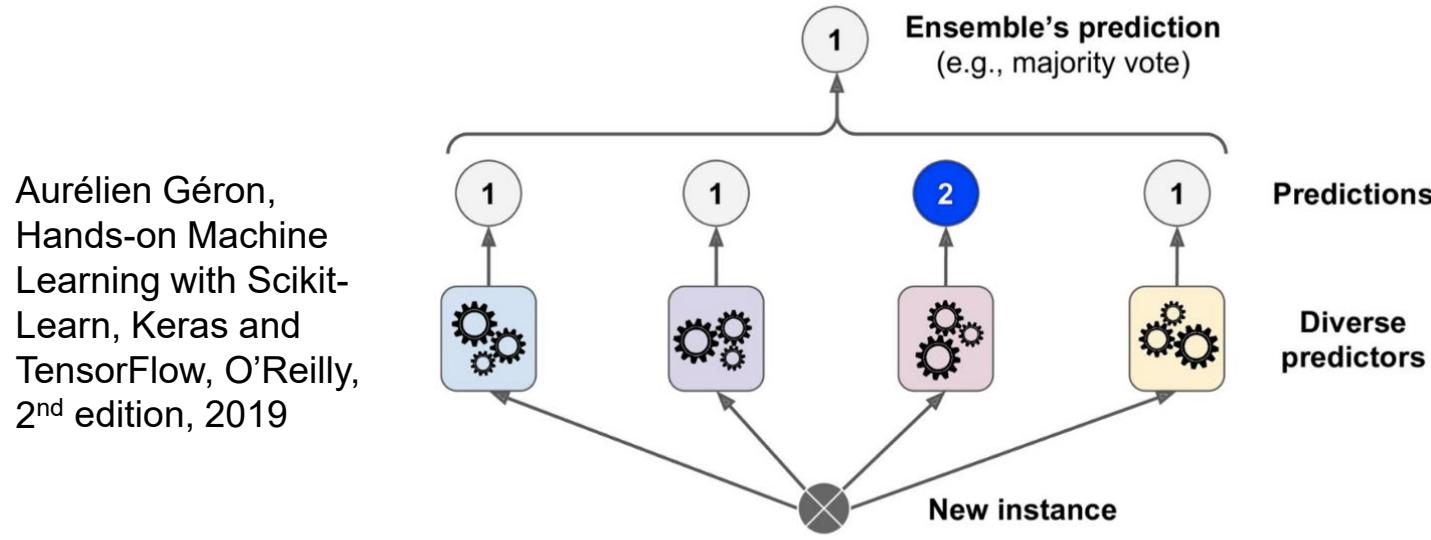
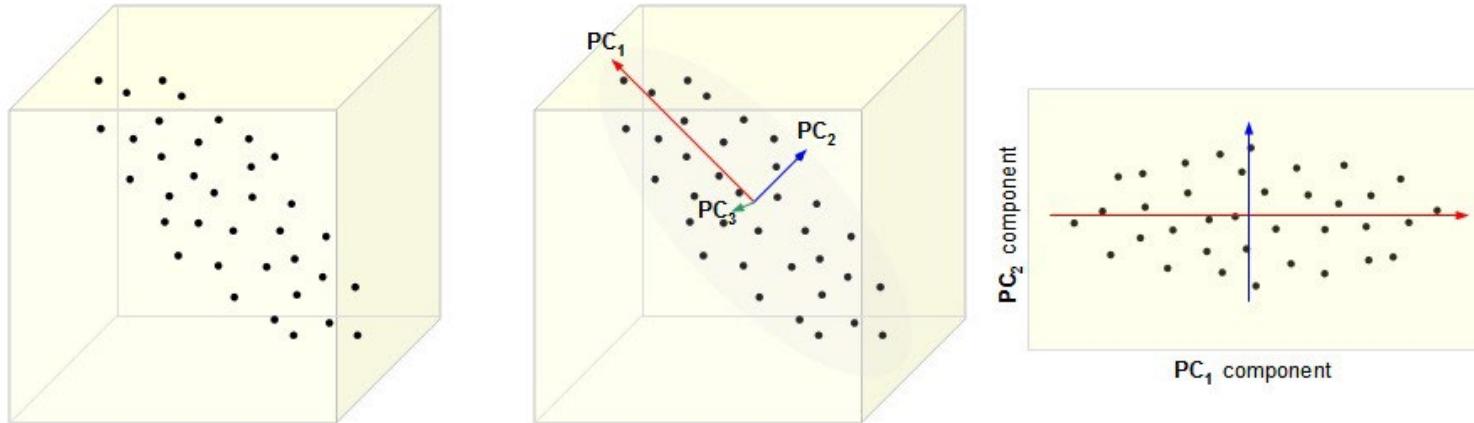


Figure 7-2. Hard voting classifier predictions

PCA (unsupervised learning)

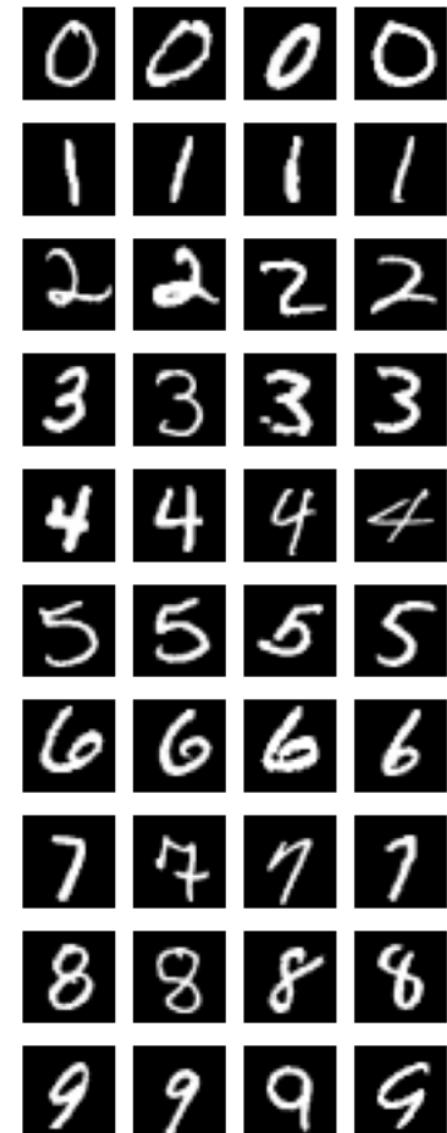
- PCA → Principal Component Analysis.
 - **Motivation:** in high dimensional data, some dimensions are redundant, others are highly correlated...
 - **Intuition:** Project the data into a low-dimensional space preserving the dimensions of maximum variance
 - SVD decomposition ($A = U\Sigma V^T$) of the training matrix $X = [x^{(1)} \dots x^{(N)}]$, the columns of V are the unit vectors of the principal components.
 - The low-d data is computed as $X_{d'} = XW_{d'}$; $W_{d'} := V[1:d']$



1-NN with and without PCA in MNIST

- Time per 10,000 queries in the test set
- Full-size input, $x \in \mathbb{R}^{784}$:
 $k=1$, accuracy=97.22%, time=849.6 s
- Removing lowest-variance dimensions with PCA, up to 98% time reduction and same accuracy!

```
dims after PCA=10, accuracy=91.95%, time=1.3 s
dims after PCA=20, accuracy=96.60%, time=6.1 s
dims after PCA=30, accuracy=97.60%, time=14.8 s
dims after PCA=40, accuracy=97.43%, time=24.6 s
dims after PCA=50, accuracy=97.45%, time=36.3 s
dims after PCA=60, accuracy=97.56%, time=48.9 s
dims after PCA=70, accuracy=97.48%, time=62.0 s
dims after PCA=80, accuracy=97.51%, time=69.3 s
dims after PCA=90, accuracy=97.44%, time=81.7 s
dims after PCA=100, accuracy=97.54%, time=93.8 s
dims after PCA=200, accuracy=97.29%, time=159.8 s
dims after PCA=300, accuracy=97.24%, time=243.2 s
dims after PCA=400, accuracy=97.22%, time=318.2 s
dims after PCA=500, accuracy=97.24%, time=400.1 s
dims after PCA=600, accuracy=97.22%, time=473.3 s
```



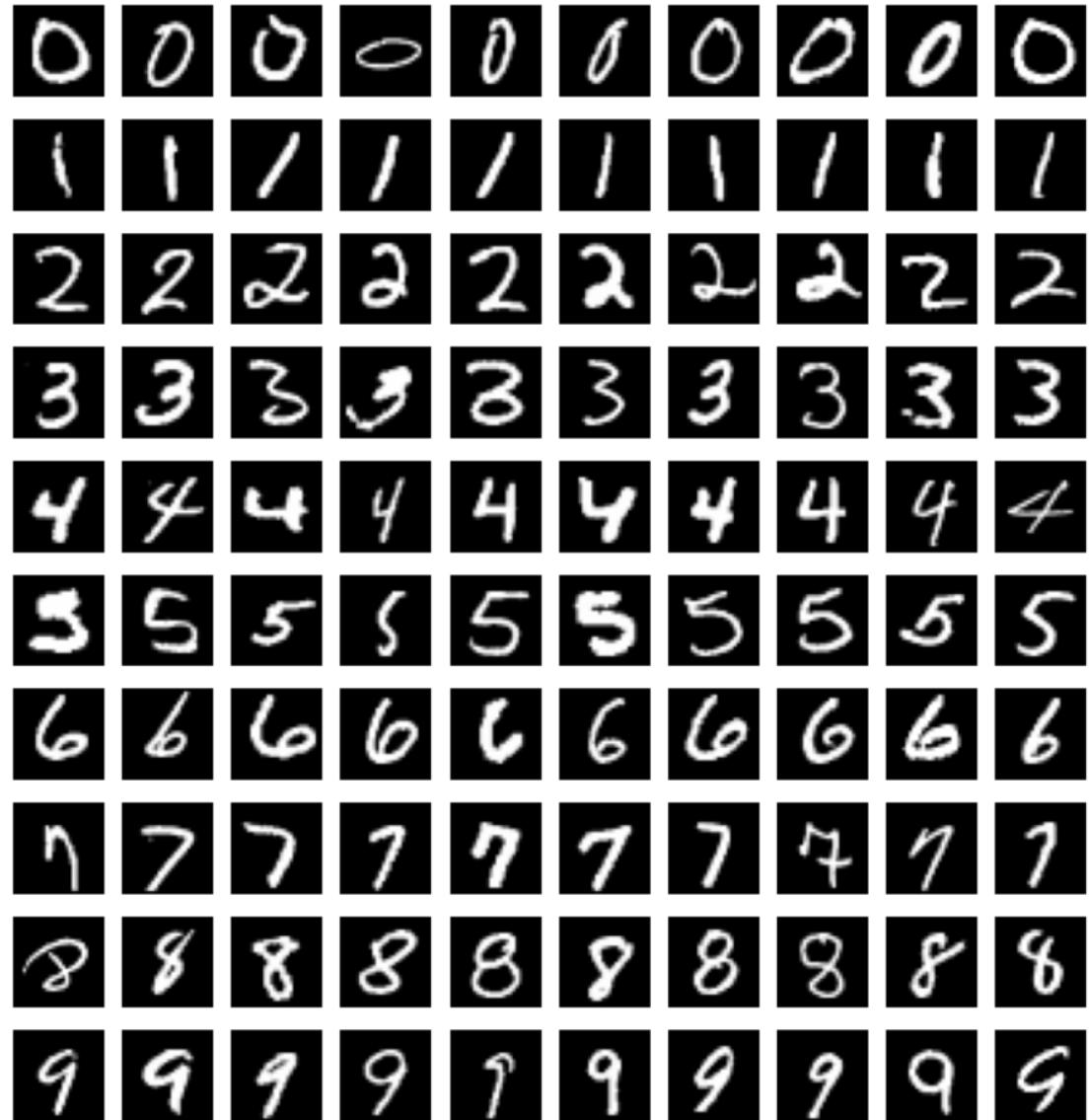
Colab link:

<https://colab.research.google.com/drive/1EdemME5vbZr01Q6CijYh0eo3hnRsjiKE>

Assignment: Logistic regression in MNIST

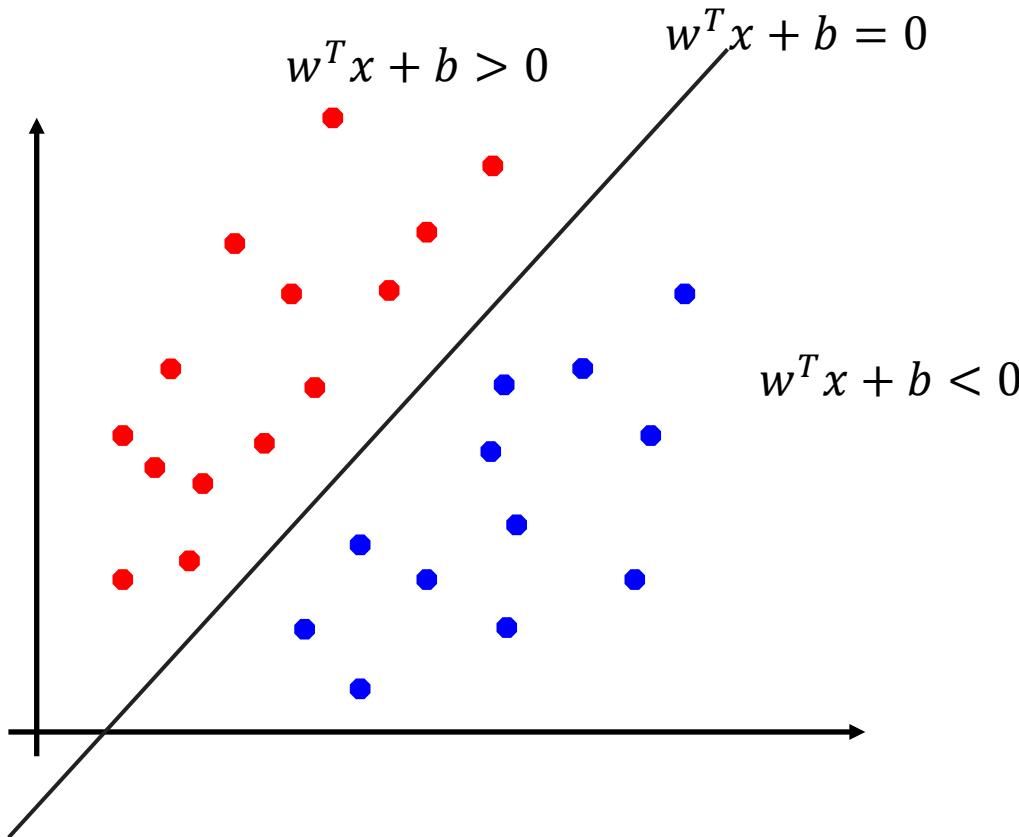
- Spoilers:
 - Excruciatingly slow!
 - Due to high dimensionality and large dataset size
 - Shortlist promising models using a reduced version of the data!
 - Use PCA!
 - 1-NN better than logistic regression...
 - why?

https://colab.research.google.com/drive/1HS9988xVWD6yqw9FRbe4z7qQfAFejL_V



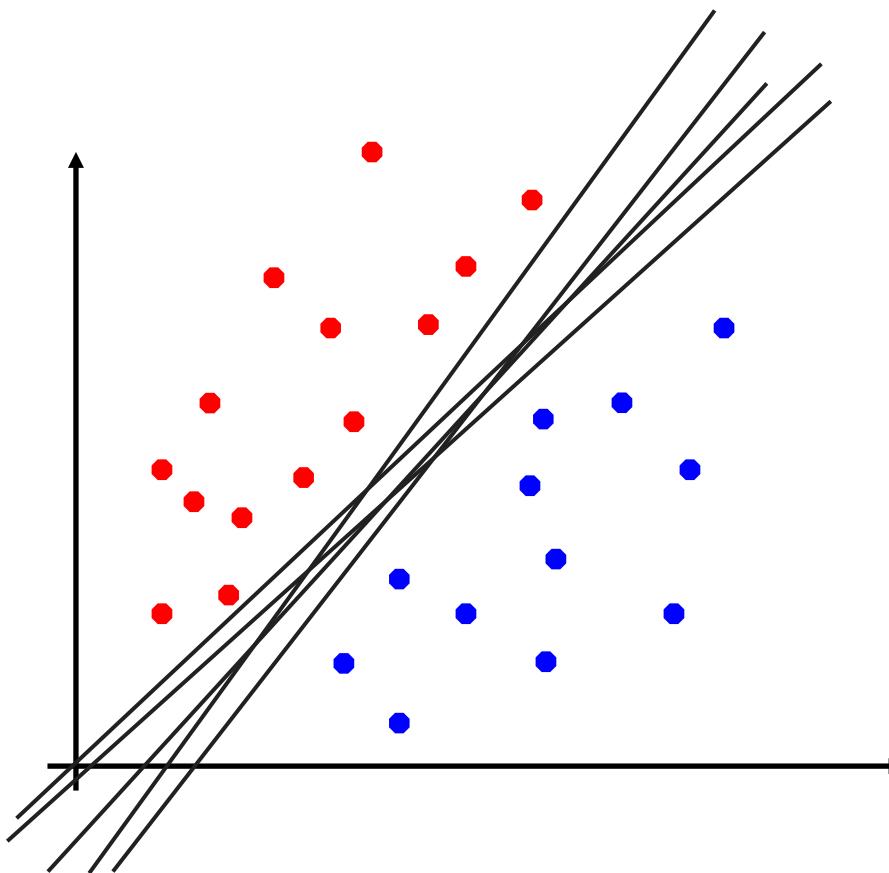
Support Vector Machines (SVM)

- Recap: Linear classification consists of finding an hyperplane in the feature space that separates the training data



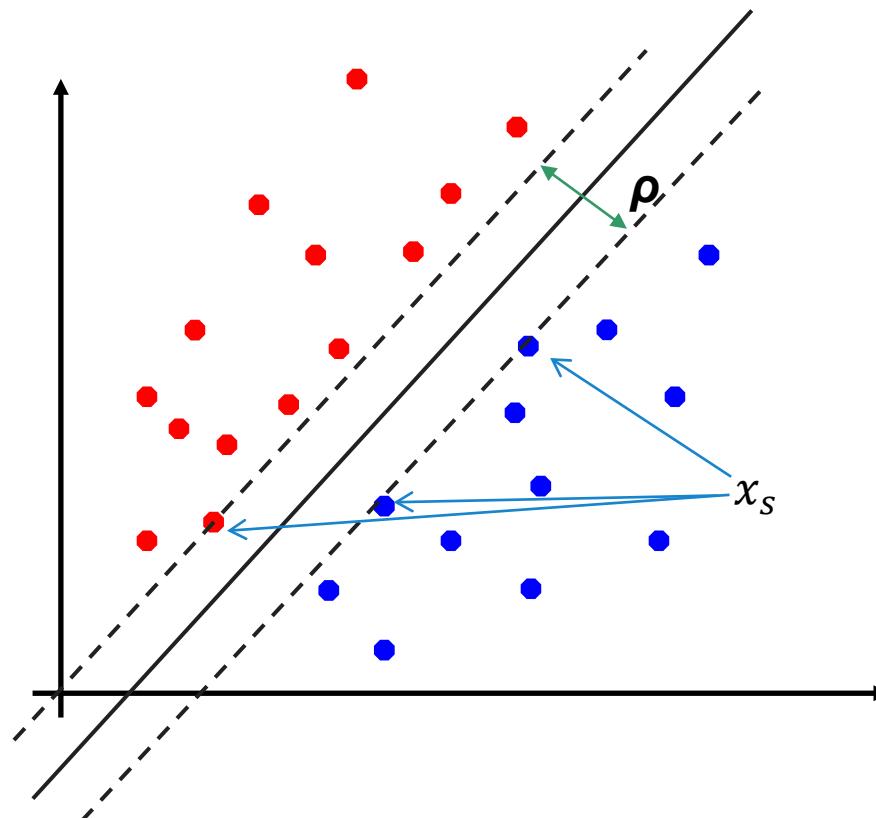
Support Vector Machines (SVM)

- But... which one of these hyperplanes would be the best one?



Support Vector Machines (SVM)

- SVM idea: let's take the one with **maximum margin** ρ
- Maximum margin:
 - The closest samples to the hyperplane are the support vectors x_s
 - The margin is the distance between the support vectors
 - Only support vectors matter! (which seems reasonable...)



“Hard-margin” SVM

- For each training example $(x^{(i)}, y^{(i)}), y^{(i)} \in \{-1, +1\}$

$$\begin{aligned} w^T x^{(i)} + b &\leq -\rho/2 \quad \text{if } y^{(i)} = -1 \\ w^T x^{(i)} + b &\geq \rho/2 \quad \text{if } y^{(i)} = +1 \end{aligned} \Rightarrow y^{(i)}(w^T x^{(i)} + b) \geq \rho/2$$

- If $x_s^{(i)}$ is a support vector, the above inequality is an equality
- After re-scaling w and b by $\rho/2$, the distance between each support vector $x_s^{(i)}$ and the hyperplane is $1/\|w\|$
- The margin is then $\rho = 2/\|w\|$
- And the margin maximization can be formulated as

$$\hat{w}, \hat{b} = \arg \max_{w,b} 2/\|w\| \text{ (alternatively } \arg \min_{w,b} \frac{1}{2} \|w\|^2)$$

subject to $y^{(i)}(w^T x^{(i)} + b) \geq 1$ (for re-scaled w and b)

“Soft-margin” SVM

- “Hard-margin SVM” does not allow data to be on the wrong side of the hyperplane
- For noisy data, we might want to allow that
- Idea: introduce a per-sample slack variable $\xi^{(i)} > 0$ modelling deviations from the right side of the hyperplane
- The optimization is now

$$\hat{w}, \hat{b} = \arg \min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_i \xi^{(i)}$$

subject to $y^{(i)}(w^T x^{(i)} + b) \geq 1 - \xi^{(i)}$ (re-scaled w and b)

- C trades off the margin and the amount of slack we have
 - $C \downarrow$ regularizes the solution (favors simple decision functions)
 - Usually tuned by cross-validation

Non-linear SVM

- The SVM parameters are learnt by solving the dual form of the optimization

$$\hat{\alpha} = \arg \max_{\alpha} \frac{1}{2} \sum_i \sum_j y^{(i)} y^{(j)} \alpha^{(i)} \alpha^{(j)} \langle \mathbf{x}^{(i)}, \mathbf{x}^{(j)} \rangle - \sum_i \alpha^{(i)}$$

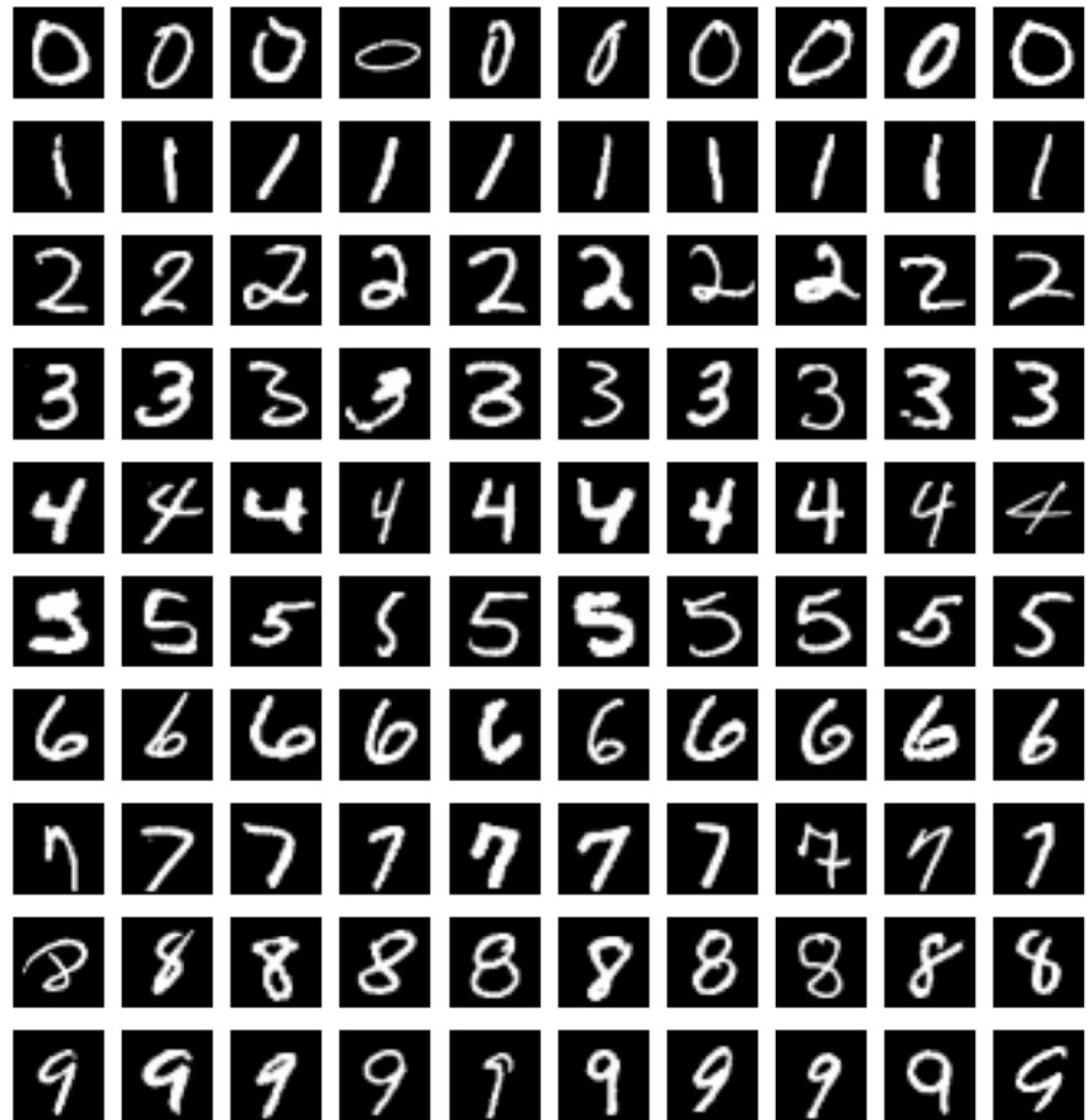
subject to $\sum_i y^{(i)} \alpha^{(i)} = 0, 0 \leq \alpha^{(i)} \leq C$

- The objective function is based in the inner product $\langle \mathbf{x}^{(i)}, \mathbf{x}^{(j)} \rangle$
- For non-linear functions of the data $\phi(x^{(i)})$, the only modification is on the inner product $\langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)}) \rangle$
- Instead of choosing the non-linear function $\phi(x^{(i)})$, we choose the similarity function $k(x^{(i)}, x^{(j)}) = \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle_{\mathcal{H}}$
- The similarity function $k(x^{(i)}, x^{(j)})$ is a kernel (more about kernels later in the course...)

Non-linear SVM hyperparameters

- Obtained by cross-validation
- C parameter
 - Trades off the margin and the slack allowed to the data
- Kernel and kernel parameters. Usual choices
 - Linear kernel: $k(x^{(i)}, x^{(j)}) = x^{(i)}x^{(j)}$
 - Polynomial kernel: $k(x^{(i)}, x^{(j)}) = (\gamma x^{(i)}x^{(j)} + r)^d$
 - Gaussian kernel: $k(x^{(i)}, x^{(j)}) = \exp\left(-\frac{\|x^{(i)} - x^{(j)}\|^2}{\sigma^2}\right)$
 - Radial basis function: $k(x^{(i)}, x^{(j)}) = \exp\left(-\gamma \|x^{(i)} - x^{(j)}\|^2\right)$
 - Sigmoid kernel: $k(x^{(i)}, x^{(j)}) = \tanh(\gamma x^{(i)}x^{(j)} + r)$

Assignment: SVM in MNIST??



https://colab.research.google.com/drive/1HS9988xVWD6yqw9FRbe4z7qQfAFejL_V

Regression metrics

- **Mean absolute error** $\rightarrow \frac{1}{N} \sum_{i=1}^N |y^{(i)} - w^T x^{(i)}|$
 - Small influence of large errors
- **Mean squared error** $\rightarrow \frac{1}{N} \sum_{i=1}^N (y^{(i)} - w^T x^{(i)})^2$
 - Highly influenced by large errors
- **Root mean squared error** $\rightarrow \sqrt{\frac{1}{N} \sum_{i=1}^N (y^{(i)} - w^T x^{(i)})^2}$
- **Median absolute error** $\rightarrow \text{median}(|y^{(1)} - w^T x^{(1)}|, \dots, |y^{(N)} - w^T x^{(N)}|)$
 - Robust to outliers!
- **R² score** $\rightarrow 1 - \frac{\sum_{i=1}^N (y^{(i)} - w^T x^{(i)})^2}{\sum_{i=1}^N (y^{(i)} - \bar{y})^2} \in [-\infty, 1]$, $\bar{y} = \frac{1}{N} \sum_{i=1}^N y^{(i)}$
 - Proportion of variance in y explained by the variables of the model
- Many others! Maximum error, Mean absolute percentage error, explained variance score, etc.
- Compare against dumb baseline (mean of the training predictions?)

Types of error (binary classification)

- **True Positive (TP):** Positive data (correctly) classified as positive data.
- **True Negative (TN):** Negative data (correctly) classified as negative data.
- **False Positive / Type I Error (FP):** Negative data (not correctly) classified as positive data.
- **False Negative / Type II Error (FN):** Positive data (not correctly) classified as negative data.

		Ground Truth	
		Positive	Negative
Prediction	Positive	TP	FP
	Negative	FN	TN

Metrics: Precision and Recall

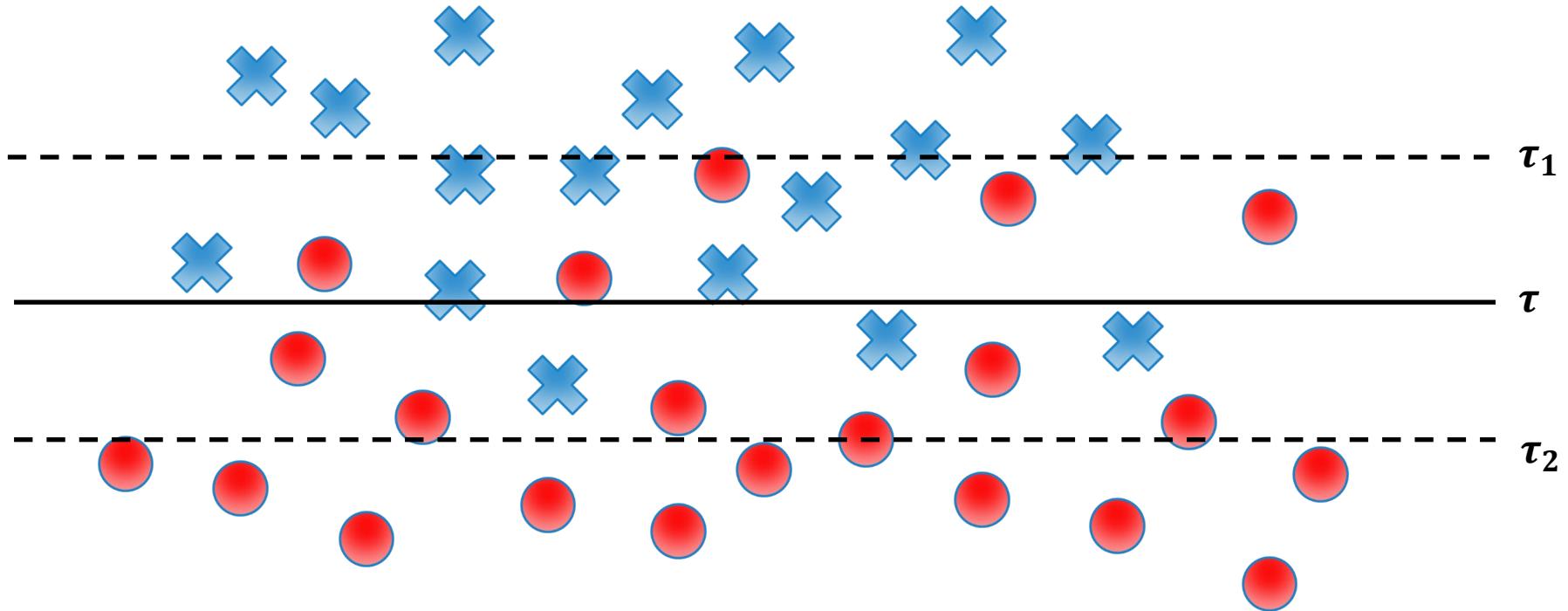
- **Accuracy, success rate, hit rate** → $\frac{TP+TN}{TP+FP+TN+FN}$
 - Limitation: Not every error is equally important...
 - Best and worst success rate?
- **Precision** → $\frac{TP}{TP+FP} \in [0,1]$
 - Ratio of the true positives over all the positive classifications.
 - Fraction of retrieved elements that are relevant.
- **Recall** → $\frac{TP}{TP+FN} \in [0,1]$
 - Ratio of the true positives over all the positive samples.
 - Fraction of relevant elements that are retrieved.
- Both depend on the classification threshold.
- They are dependent, if one increases the other decreases. Sometimes given in pairs (e.g., precision at recall X).

Precision and Recall (some intuition)

- The girlfriend (boyfriend) example:
 - Your girlfriend/boyfriend makes birthday presents for you every year.
 - One year she/he asks: Do you remember all my presents from the last 10 years?
 - **Precision** is how many are correct from the ones you said you remember.
 - **Recall** is how many correct you can remember from all correct ones.
 - How can we have 100% precision or recall? What happens to the other?
 - What is the relation to the classification threshold?
- The airport security example:
 - FP and FN: Have both errors the same consequences?
 - How can we have 100% precision or recall?
- Do you need precision, recall or both?
 - The solution is, most often, a compromise...

Precision-Recall and thresholds

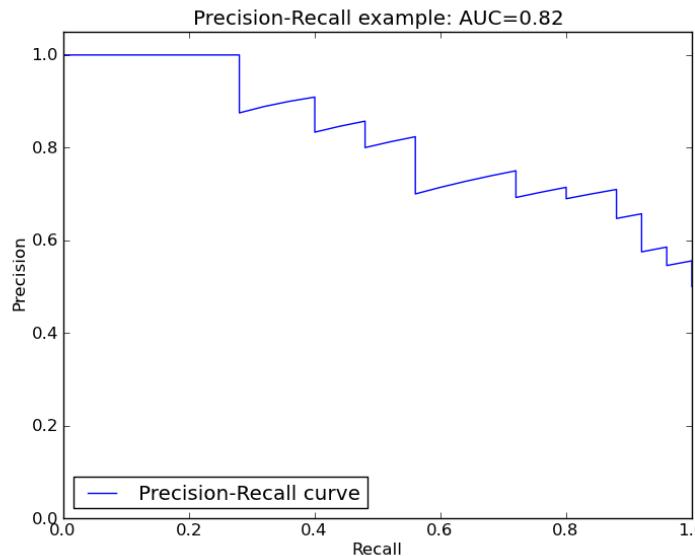
- Two-class problem



- What happens to the TP, TN, FP, FN, precision and recall if we move the boundary from τ to τ_1 or τ_2 ?

Precision-Recall curve

- 2D plot, recall in x-axis and precision in y-axis.
- It shows the precision-recall pairs when we change the classifier threshold τ
- $\tau \uparrow \Rightarrow FP \downarrow \Rightarrow FN \uparrow \Rightarrow P \rightarrow 1 \Rightarrow R \rightarrow 0$
- $\tau \downarrow \Rightarrow FP \uparrow \Rightarrow FN \downarrow \Rightarrow P \rightarrow 0 \Rightarrow R \rightarrow 1$
- Changing a classifier threshold τ it operates in different modes. We have to select the more appropriate for our purposes.
- For general purposes, the area under the curve (AUC) is usually taken as a general metric (the higher the better).
- What curve and AUC would have a perfect classifier? What curve and AUC would have the worst possible classifier?



F1 score, confusion matrix

F1 score

- Harmonic mean of precision and recall (the higher the better)

$$F_1 = \frac{2}{\frac{1}{P} + \frac{1}{R}} = 2 \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}} \in [0,1]$$

- It can be weighed to give β times more importance to precision than recall $F_1 = (1 + \beta^2) \frac{\text{precision} * \text{recall}}{\beta^2 * \text{precision} + \text{recall}}$

Confusion matrix

- Rows: true categories.
- Columns: predicted categories.
- Cells: number/frequency of occurrence.
- More information than other metrics (for example, typical classification mistakes)
- Difficult to display if many categories.
- Difficult to summarize...

	Dive	.86	.00	.00	.00	.00	.07	.00	.00	.00	.07
Dive	.86	.00	.00	.00	.00	.00	.07	.00	.00	.00	.07
Golf	.00	.94	.00	.00	.00	.06	.00	.00	.00	.00	.00
Kick	.00	.00	.75	.00	.05	.15	.00	.00	.00	.00	.05
W.Lift	.00	.00	.00	1.0	.00	.00	.00	.00	.00	.00	.00
Ride	.00	.00	.08	.00	.92	.00	.00	.00	.00	.00	.00
Run	.00	.08	.38	.00	.08	.46	.00	.00	.00	.00	.00
SK.Board	.00	.00	.08	.00	.08	.00	.83	.00	.00	.00	.00
Swing 1	.00	.00	.00	.00	.00	.00	.00	1.0	.00	.00	.00
Swing 2	.00	.00	.00	.00	.00	.00	.00	.00	1.0	.00	.00
Walk	.00	.23	.05	.00	.05	.05	.05	.05	.00	.00	.59

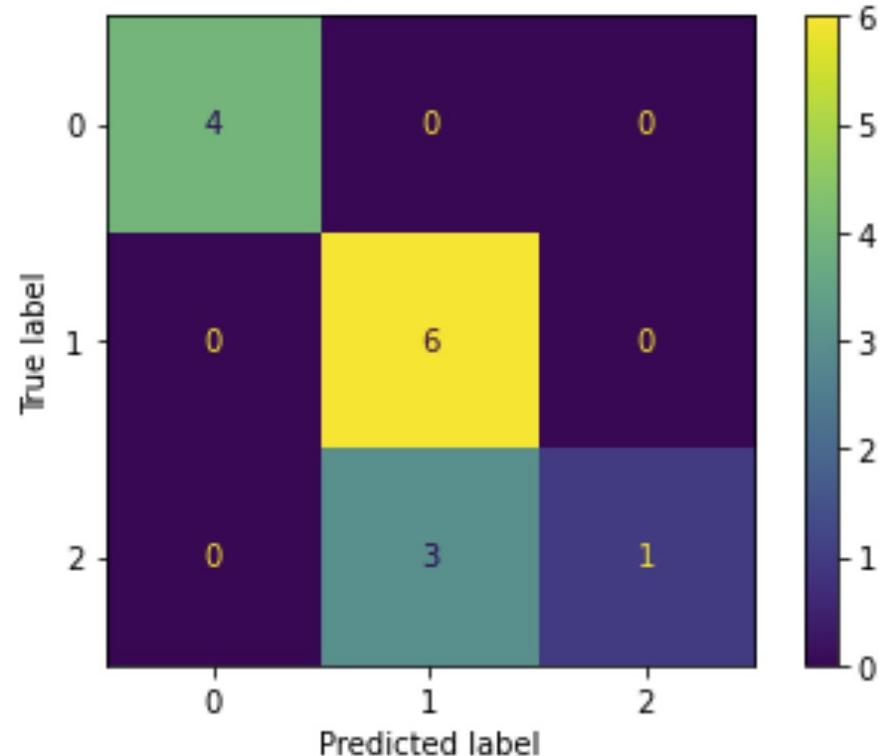
Dive Golf Kick W.Lift Ride Run SK.Board Swing 1 Swing 2 Walk

Confusion matrix for action recognition from videos

Taken from Qiu, Jiang, and Chellappa, "Sparse Dictionary-based Representation and Recognition of Human Action Attributes", ICCV 2011

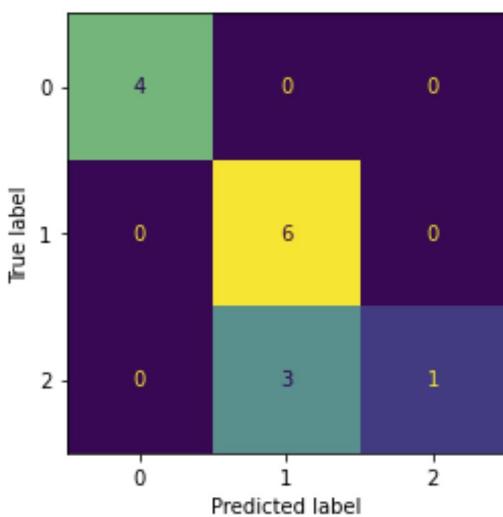
TP, TN, FP and FN from the confusion matrix

- Remember:
 - Rows: true categories.
 - Columns: predicted categories.
- **True positives per class:**
 - Diagonal elements [4, 6, 1]
- **False positives per class:**
 - Sum of the off-diagonal elements of each column [0, 3, 0]
- **False negatives per class:**
 - Sum of the off-diagonal elements of each row [0, 0, 3]
- **True negatives per class:**
 - Sum of all matrix elements except the row and column of the class [10, 5, 10]



Multiclass Precision/Recall

- True positives per class: [4, 6, 1], True negatives per class: [10, 5, 10]
- False positives per class: [0, 3, 0], False negatives per class: [0, 0, 3]



	Accuracy	Precision	Recall
Class 1	$\frac{4+10}{4+6+1+3} = 1.0$	$\frac{4}{4+0} = 1.0$	$\frac{4}{4+0} = 1.0$
Class 2	$\frac{6+5}{4+6+1+3} = 0.79$	$\frac{6}{6+3} = 0.66$	$\frac{6}{6+0} = 1.0$
Class 3	$\frac{1+10}{4+6+1+3} = 0.79$	$\frac{1}{1+0} = 1.0$	$\frac{1}{1+3} = 0.25$

- Overall accuracy = $\frac{\text{correct}}{\text{total}} = \frac{4+6+1}{4+6+1+3} = 0.78$
- Average accuracy = $\frac{1.0+0.79+0.79}{3} = 0.86$
- Average Precision = $\frac{1.0+0.66+1.0}{3} = 0.89$
- Average Recall = $\frac{1.0+1.0+0.25}{3} = 0.75$

If data is imbalanced, average metrics may be weighted...

Machine Learning (69152)

Máster in Robotics, Graphics
and Computer Vision

