

Lecture 03 30.09.2024

Regularized regression

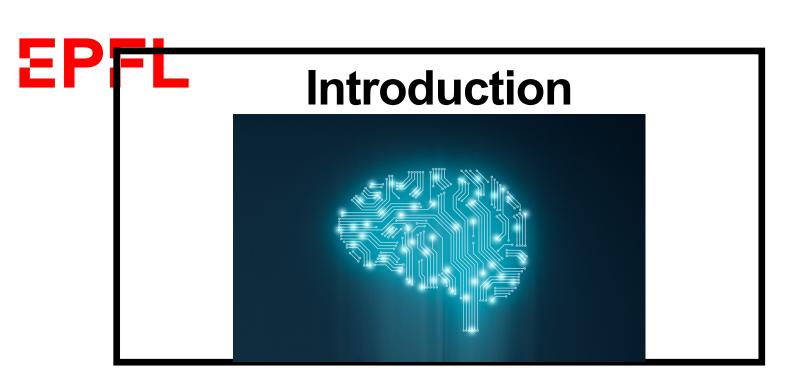
Classification through logistic regression

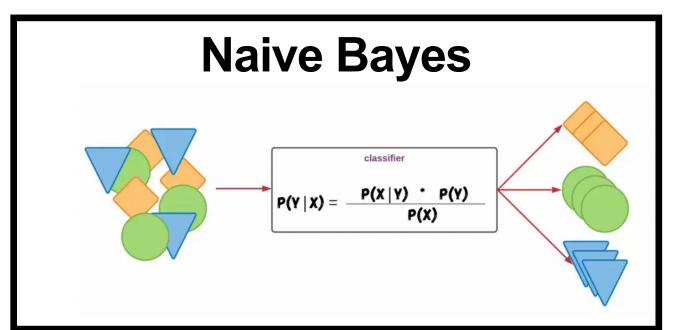


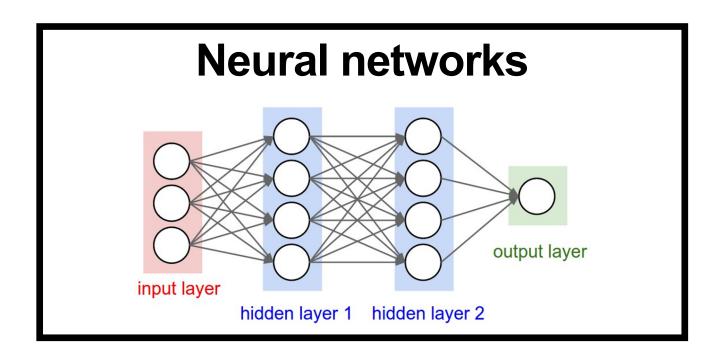
Outline

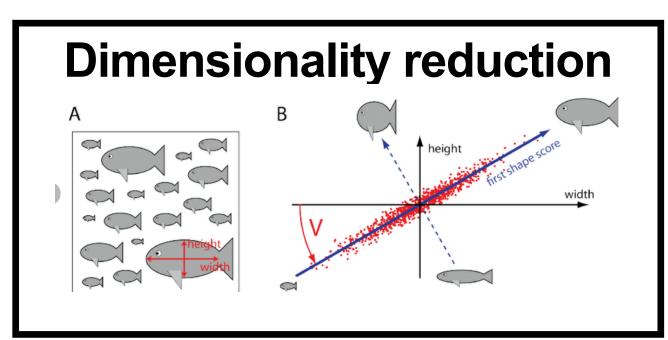
- Overfitting/underfitting and regularization
- Logistic regression

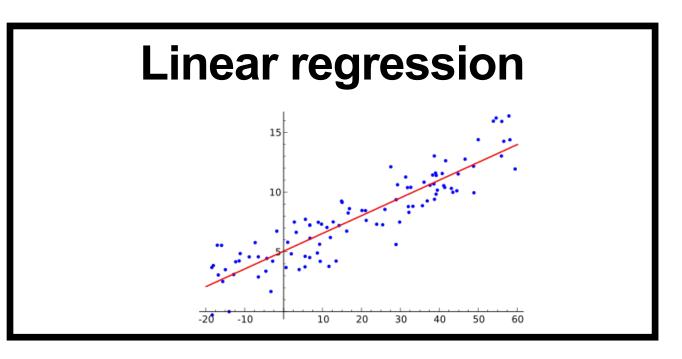
- Announcements:
 - Check moodle for past year's exams and quizzes

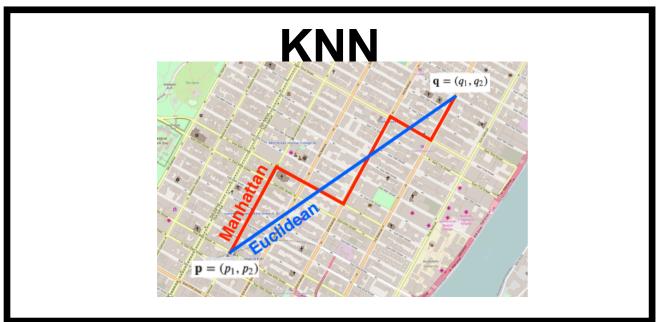


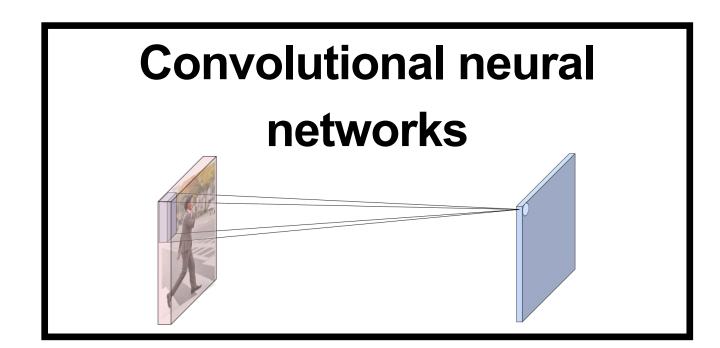


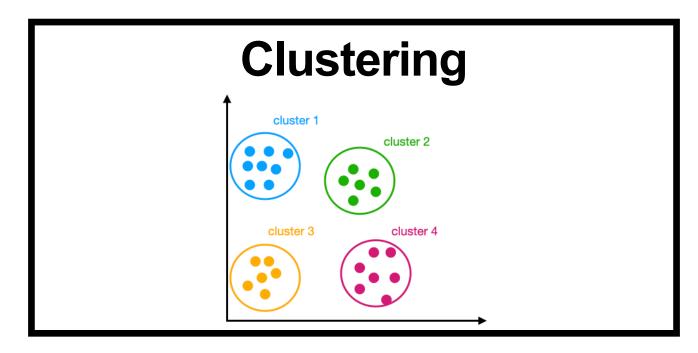


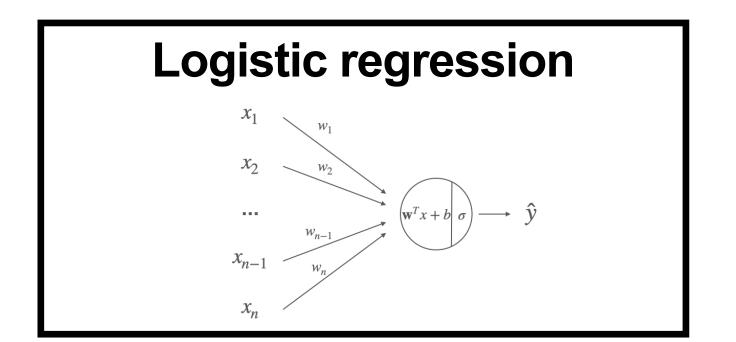


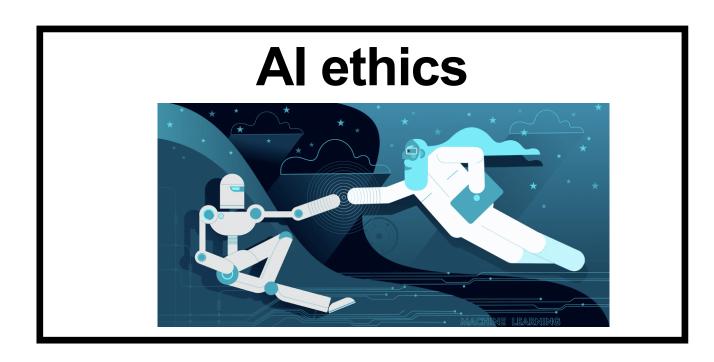


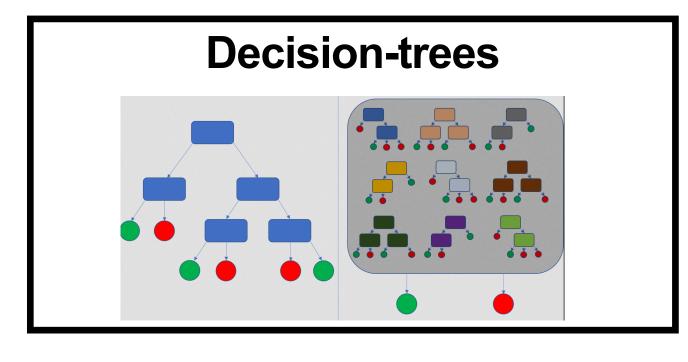


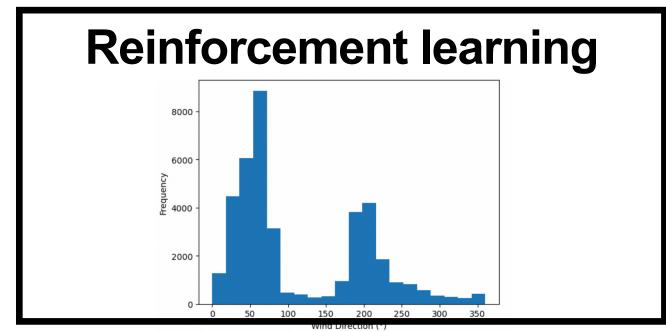














Review of linear regression

• Linear model:
$$\{x', y'\}_{i=1}^{N}$$
, $x' \in \mathbb{R}^{N}$. $y' \in \mathbb{R}^{N}$
 $y = b + \omega_{1} \times 1 + \cdots + \omega_{d} \times d$ $\omega = [b, \omega_{1}]$
 $\psi : \mathbb{R}^{d} \rightarrow \mathbb{R}^{p}$
 $y = b + \omega_{1} \oplus (x_{1}) + \cdots + \omega_{p} \oplus_{p} (x_{p})$

Mean-Square-Error (MSE) loss

$$\dot{y}$$
 = $b + w, x, + \dots + w, x, d$

$$J(\omega) = \frac{1}{N} \left(\frac{3}{3} - \frac{1}{3} \right)^{2}$$

 $w = [b, w_1, ..., w_d] \in \mathbb{R}$



Underfitting and overfitting



Underfitting & Overfitting

Goal of supervised ML models: generalise well on new data (based on the patterns learned from known data).

Ways ML can fail

- · Underfitting = model too simple, not enough doctor
- Overfitting





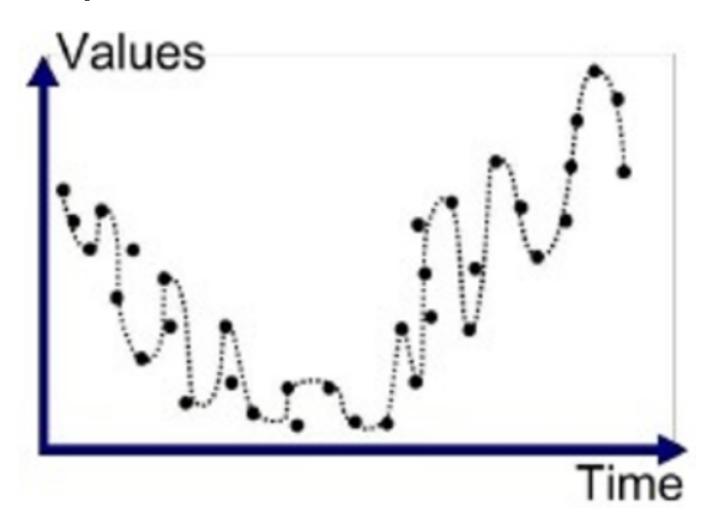
Underfitting & Overfitting Overfitting

Q: What if the ML model is too complex?

Overfitting model:

- Fits well on training data
- Doesn't generalise well to unknown data.

Reason: the model is too complex \rightarrow it fits the noises and errors.



The model passes by every data point but doesn't capture the U shape of the data set.

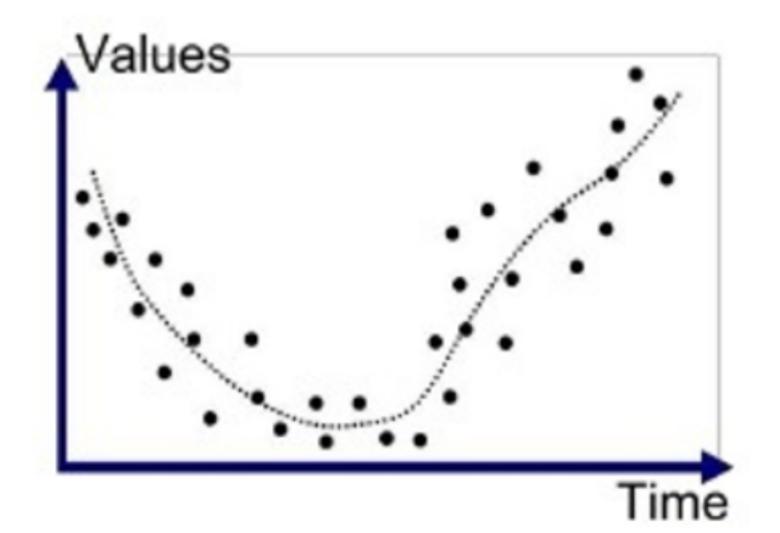


Underfitting & Overfitting Overfitting - Solutions

Solutions:

- Simpler model → fit the data and not the noises and errors.
- More training data (→ less effect of noise)
- Add a regularisation term (common solution)

Goal: Find the equilibrium between fitting the training data and keeping the model simple enough to ensure it will generalise well on new data.





Sensitivity and Regularization

- Sensitivity of a predictor: similar data should lead to similar outcome
 - Less sensitive models tend to not overfit ∞ , $\infty' \in \mathbb{R}^{c'}$

Similar data
$$x, x'$$
 $\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{2}^{2})^{2} + \cdots + (x_{d} - x_{d}^{2})^{2}\right)^{2}$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{2}^{2})^{2} + \cdots + (x_{d} - x_{d}^{2})^{2}\right)^{2}$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{2}^{2})^{2} + \cdots + (x_{d} - x_{d}^{2})^{2}\right)^{2}$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{2}^{2})^{2} + \cdots + (x_{d} - x_{d}^{2})^{2}\right)^{2}$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{2}^{2})^{2} + \cdots + (x_{d} - x_{d}^{2})^{2}\right)^{2}$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{2}^{2})^{2} + \cdots + (x_{d} - x_{d}^{2})^{2}\right)^{2}$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{2}^{2})^{2} + \cdots + (x_{d} - x_{d}^{2})^{2}\right)^{2}$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{2}^{2})^{2} + \cdots + (x_{d} - x_{d}^{2})^{2}\right)$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{2}^{2})^{2} + \cdots + (x_{d} - x_{d}^{2})^{2}\right)$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{2}^{2})^{2} + \cdots + (x_{d} - x_{d}^{2})^{2}\right)$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{2}^{2})^{2} + \cdots + (x_{d} - x_{d}^{2})^{2}\right)$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{2}^{2})^{2} + \cdots + (x_{d} - x_{d}^{2})^{2}\right)$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{d}^{2})^{2}\right)$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{d}^{2})^{2}\right)$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{d}^{2})^{2}\right)$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{d}^{2})^{2}\right)$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{d}^{2})^{2}\right)$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{d}^{2})^{2}\right)$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{d}^{2})^{2}\right)$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{d}^{2})^{2}\right)$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{d}^{2})^{2}\right)$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{d}^{2})^{2}\right)$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{d}^{2})^{2}\right)$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left(x, -x^{2} + (x_{2} - x_{d}^{2})^{2}\right)$$

$$\left(\left\|x - x'\right\|_{2}^{2}\right) = \left($$



Sensitivity of linear models

Motivation for regularization through Lipschitz constant of the predictor

Consider
$$f_{\omega}(x) = b + w_{1}x_{1} + ... + w_{d}x_{d}$$

$$|f_{\omega}(x) - f_{\omega}(x')| = |w_{1}(x_{1} - x_{1}') + w_{2}(x_{2} - x_{2}') + ... + w_{d}(x_{d} - x_{d}')|$$

$$= ||w_{1}||^{2} ||x_{1} - x_{1}'|| ||w_{1}||^{2} ||x_{1} - x_{2}'||$$

$$= ||w_{1}||^{2} ||x_{1} - x_{1}'|| ||w_{1}||^{2} ||x_{1} - x_{2}'||$$

Cauchy- Schwarz inequality

$$|f_{\omega}(x')|| ||w_{1}||^{2} ||x_{1} - x_{2}'||$$

$$|f_{\omega}(x')|| ||f_{\omega}(x)|| ||f_{\omega}(x)||$$

$$|f_{\omega}(x')|| ||f_{\omega}(x')||$$

$$|f_{\omega}(x')|| ||f_{\omega}(x')||$$

$$|f_{\omega}(x')|| ||f_{\omega}(x')||$$

$$|f_{\omega}(x')|| ||f_{\omega}(x')||$$



Exercise - optimize a regularized loss function

• Regularised loss function
$$J(w) = \frac{1}{N} \sum_{i=1}^{N} (w^T x^i - y^i)^{2_1} + \lambda (w_1^2 + w_2^2 + \dots + w_d^2)$$

$$\lambda \in \mathbb{R}_+$$

Find the gradient of the regularized loss function with respect to the parameters

pærameters
$$w = \begin{bmatrix} b \\ w_1 \end{bmatrix} \in \mathbb{R}^{d+1}$$

• Next - how do we set the regularization parameter λ ?



Train, validate, test in ML

Setting hyperparameters: example, the regulariser

Your Dataset



Train, validate, test in ML

Setting hyperparameters: example, the regulariser

Your Dataset

Idea #1: Split data into train and test, choose hyperparameters that work best on test data

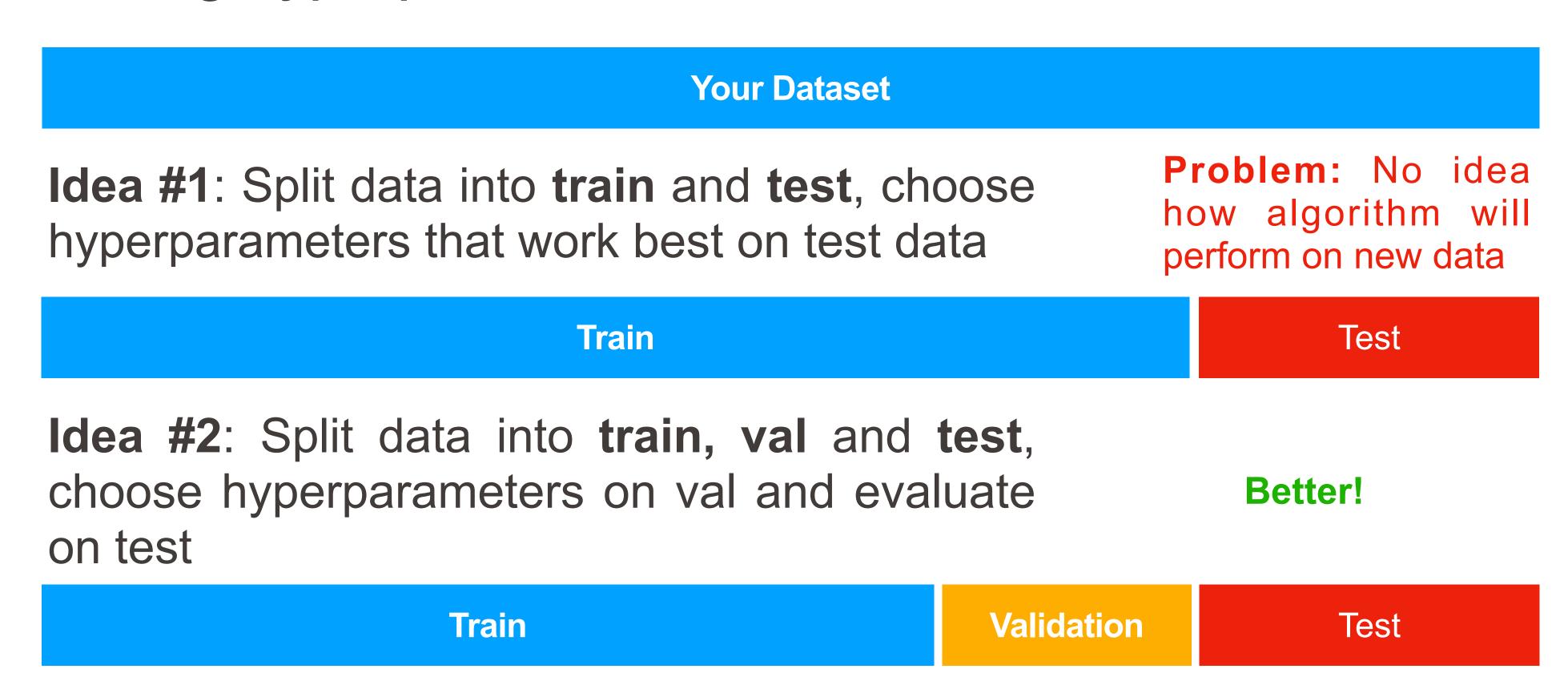
BAD: No idea how algorithm will perform on new data

Train

Test



Train, validate, test in ML Setting hyperparameters





Train, validate, test in ML

Setting hyperparameters

Your Dataset

Idea #3: Cross-Validation: Split data into folds, try each fold as validation and average the results

Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Test
Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Test
Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Test
Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Test
Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Test

Useful for small datasets

YEL

| X', y' | N Supervised unsopervised reinforcement leave leave leave leave greater leave from y' EIR y' E [1,2,...K]

Logistic Regression

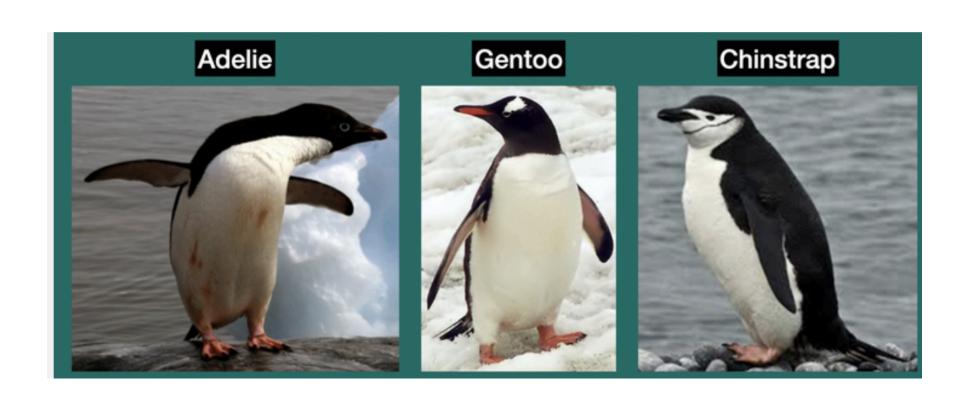
is a classification technique.

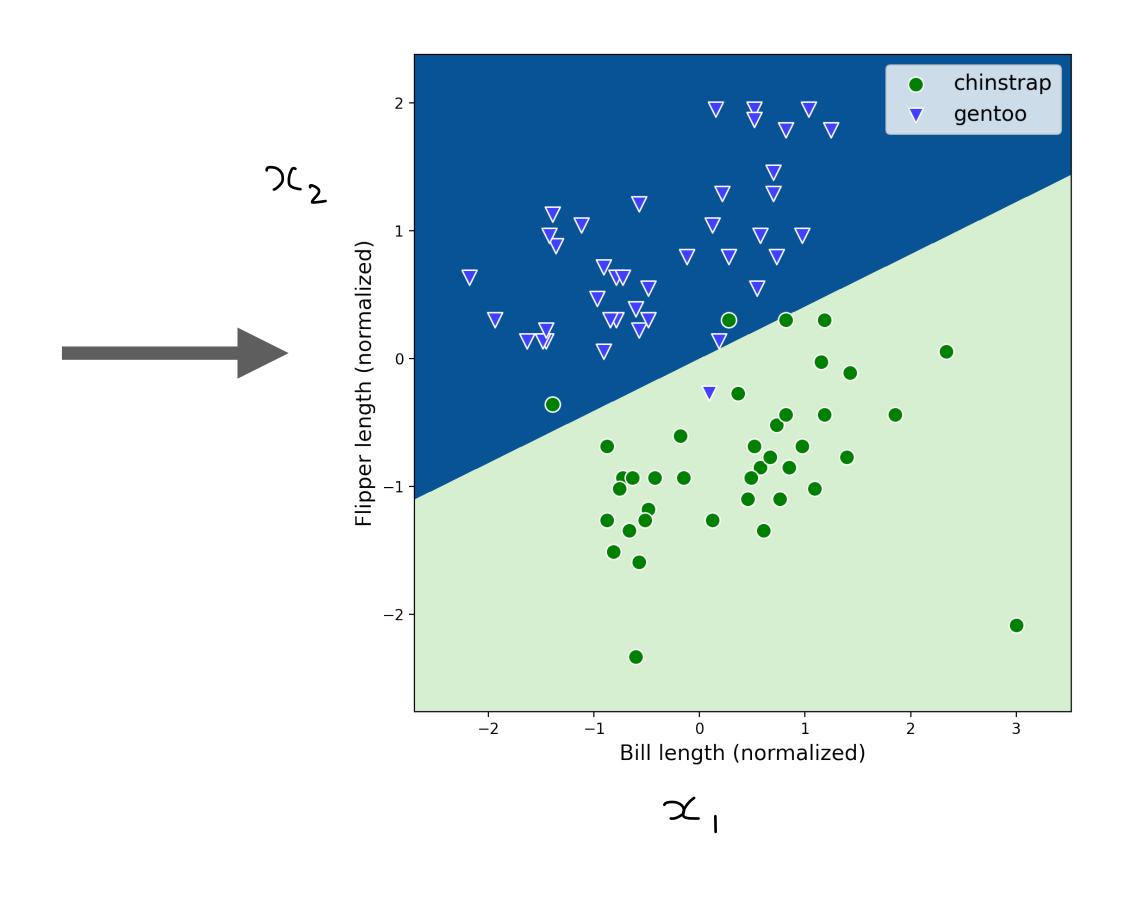


Classification

Palmer Penguins

	X (7 2	χ	7 4
	species	bill_length_mm	bill_depth_mm	flipper_length_mm	body_mass_g
0	Chinstrap	49.0	19.5	210.0	3950.0
1	Chinstrap	50.9	19.1	196.0	3550.0
2	Gentoo	42.7	13.7	208.0	3950.0
3	Chinstrap	43.5	18.1	202.0	3400.0
4	Chinstrap	49.8	17.3	198.0	3675.0



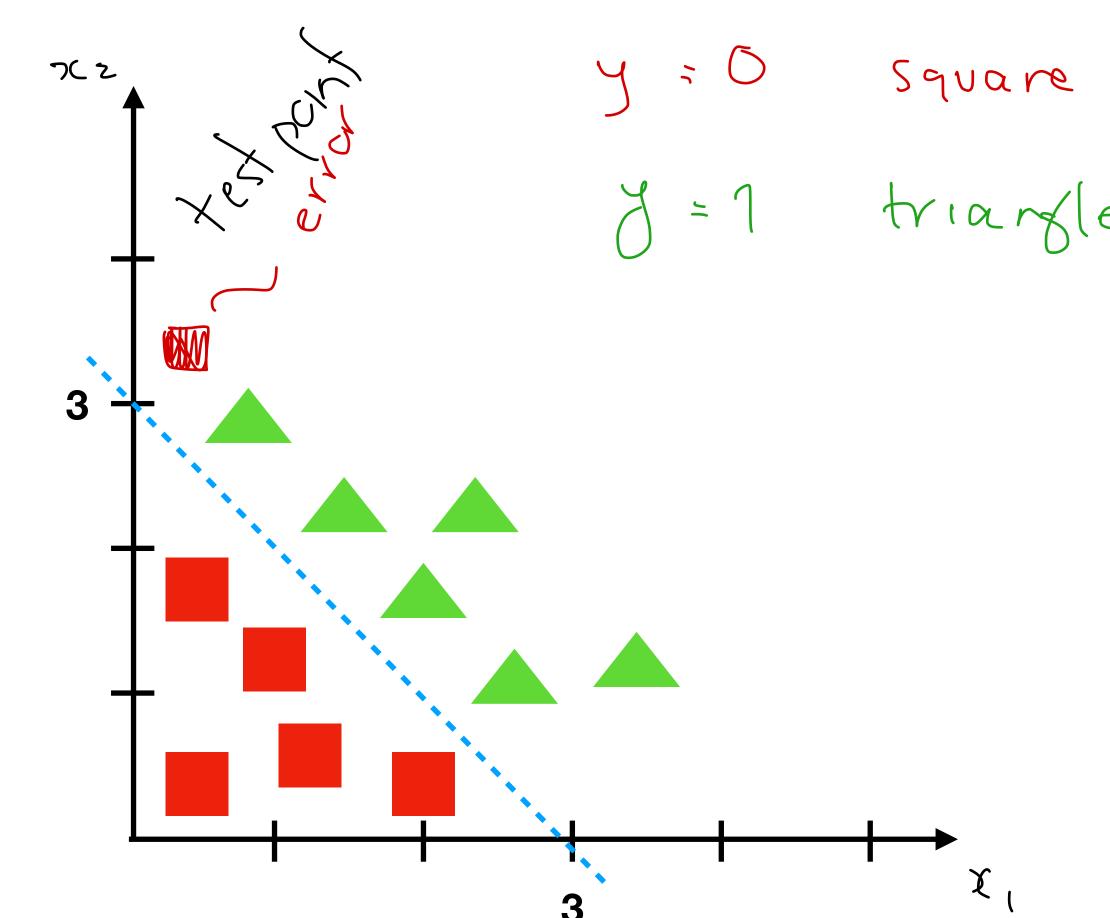




Logistic Regression (binary classification)

Characterize a classifier based on a linear model

Goal: find a line/hyperplane separating the classes



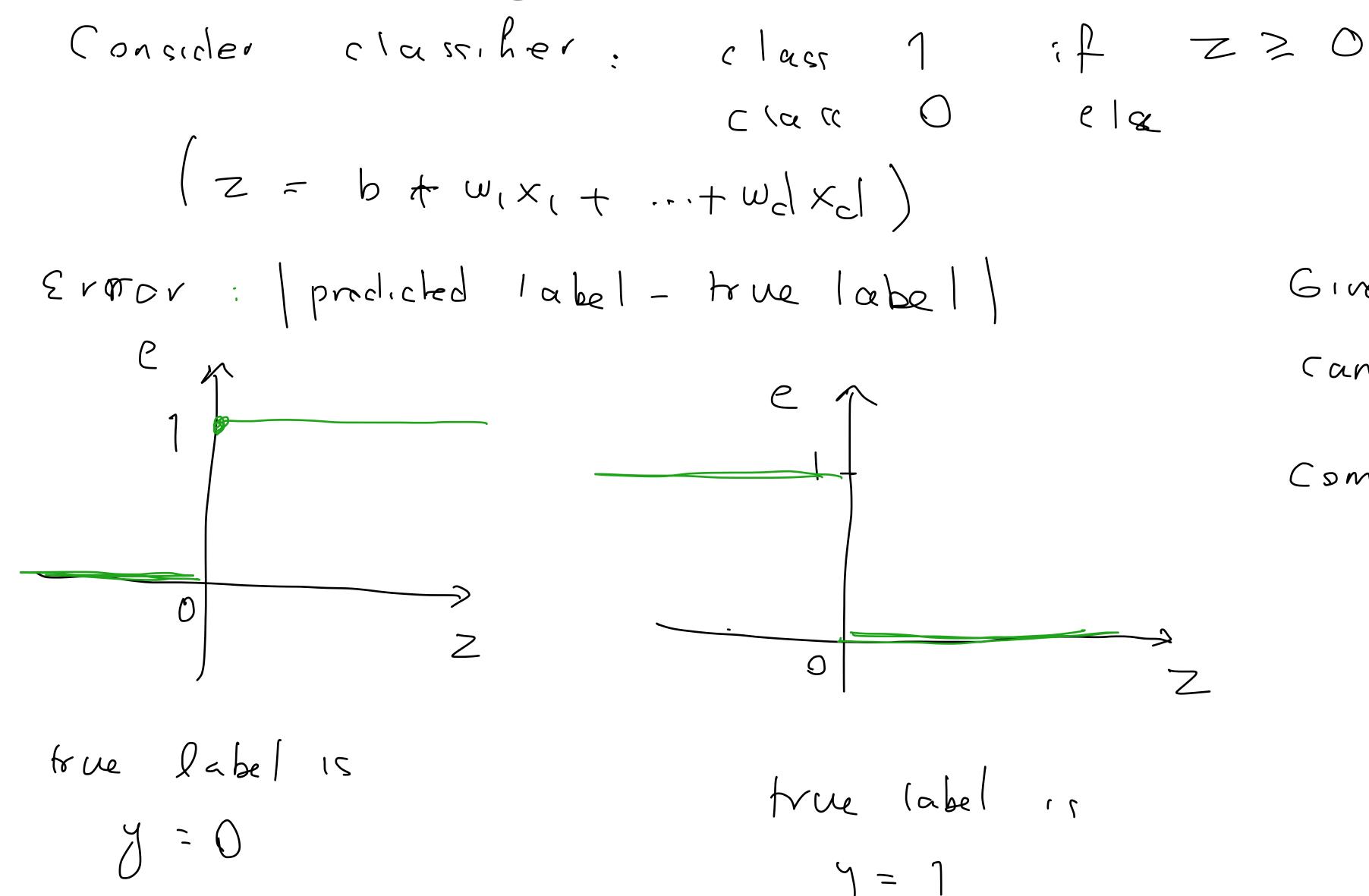
$$\mathbf{Ex: w} = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} \mathbf{b} \\ \mathbf{w}_{1} \end{bmatrix}$$

Predict triangle if
$$-3 + x_1 + x_2 \ge 0$$

$$\text{Product square else}$$



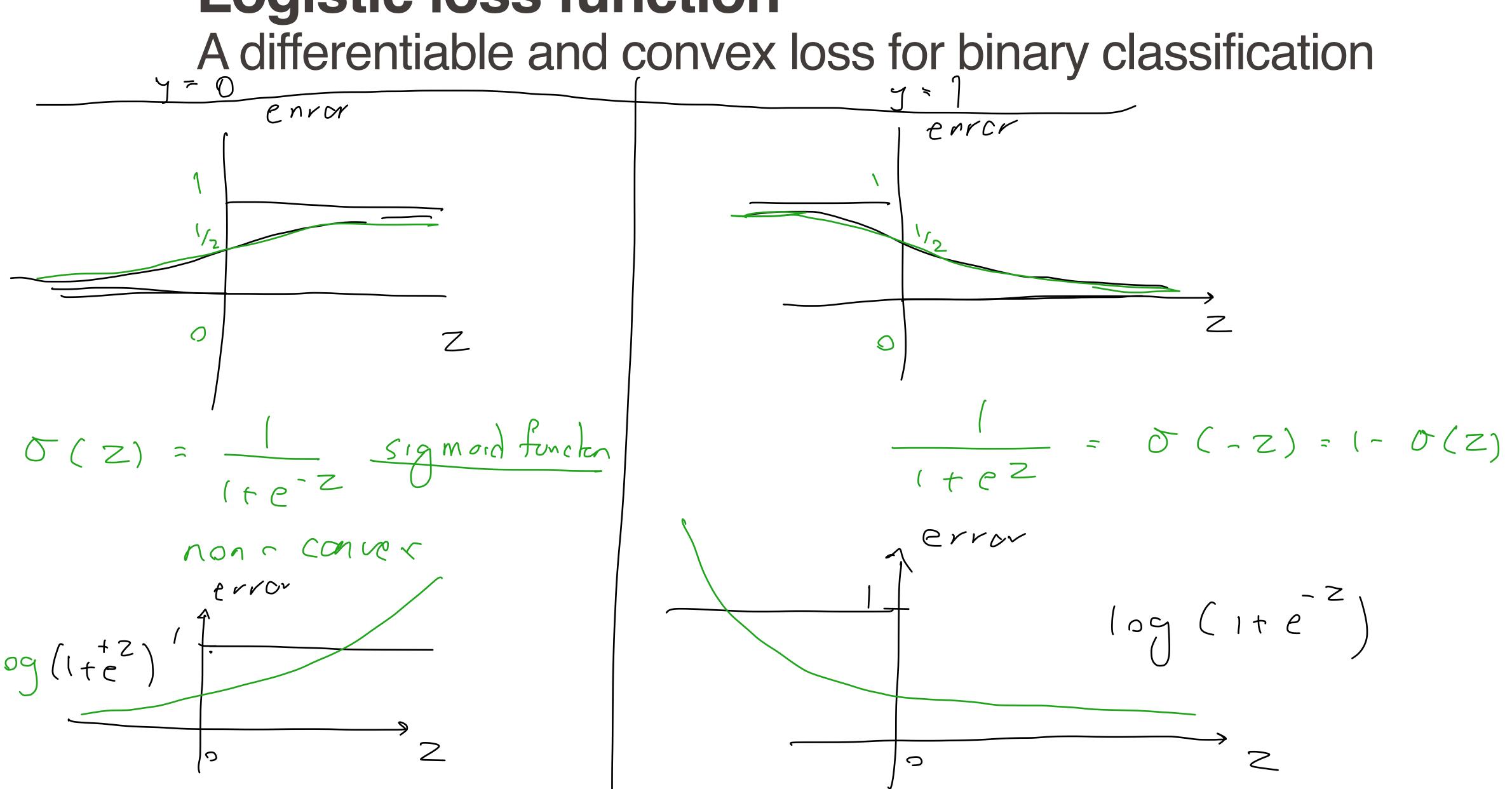
Defining loss function for classification



Given x', we can find j' &



Logistic loss function



Read more about logistic function/verify your gradient computation here: https://en.wikipedia.org/wiki/Logistic_function

EPFL

Training - Minimizing the logistic loss function

$$(1+1) \omega = (1+1) \omega = (1+1) \omega$$

EPFL

Exercise

Consider the sigmoid function $\sigma: \mathbb{R} \to (0,1), \ \sigma(z) = \frac{1}{1+e^{-z}}.$ Compute $\frac{d\sigma(z)}{dz}.$ Use the chain rule to compute $\frac{d\sigma(z(w))}{dz}$, where $z=w_0+w_1x_1+\ldots+w_dx_d$ and $w=(w_0,w_1,\ldots,w_d)$

Compute the gradient of the binary cross-entropy loss function with respect to the parameters $w = (w_0, w_1, ..., w_d)$

You can check your answers with the notes in the python exercises this week.



Probabilistic interpretation of logistic function

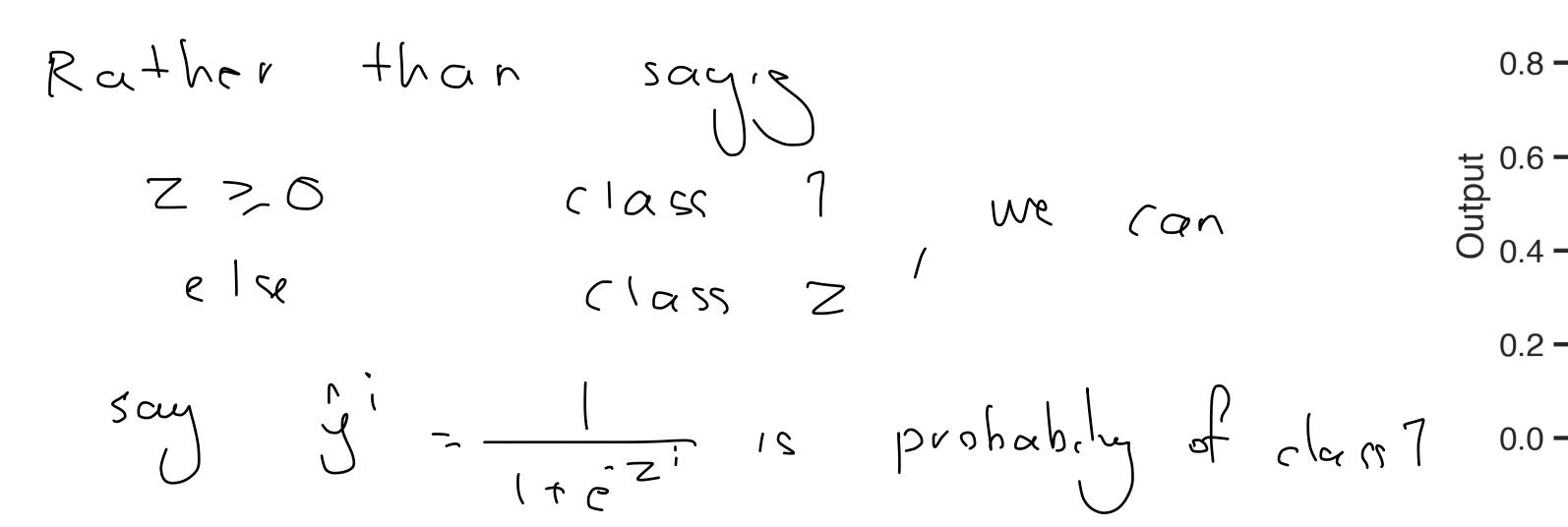
Consider the regression

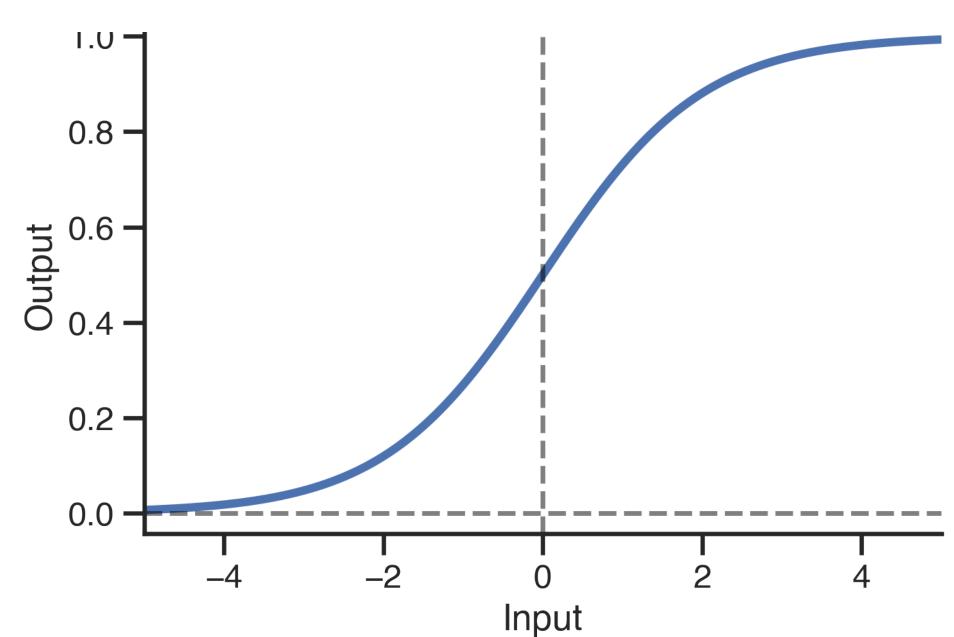
•
$$z^i = \mathbf{w}^T \mathbf{x}^i + b$$

Followed by the logistic function
$$\hat{y}^i = \frac{1}{1 + e^{-(\mathbf{w}^T x^i + b)}}$$

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

$$\sigma: \mathbb{R} \to (0,1)$$





- \hat{y}^i = estimated probability of class 1 on input x^i
- $\hat{y}^i = P(y^i = 1 \mid x^i; \mathbf{w})$ estimated probably of class 1 given x^i , parameterized by \mathbf{w}
- $-1 \hat{y}^i = P(y^i = 0 | x^i; \mathbf{w})$



Probability distributions background

$$S = \{1, 2, ..., n\}$$
, a probability dishibuter on S_{ii} a function $p: S \rightarrow \mathbb{R}$ sit. $P(i) \geq 0$ $\forall i \in S$ and $\widehat{\mathbb{Z}}_{p(i)} = 1$, we can equivently represent this function as a vector $p = \begin{bmatrix} P_1 \\ P_2 \end{bmatrix}$, where $p_i = p(i)$



Cross entropy of a distribution relative to another

probabil. Ly Let p,q: {1,2,...n} -> 1R two probability dishbuter Then, the cross entropy of 9 with respect to dishbuten p 15 defined as — $\sum_{i=1}^{n} p(i) \log (q(i))$ Our eshmated dishibuter is prohability of prediction labely given x. true dishibution is p - [p(y=1|x)(og(p(y=1|x))+]Cross. enhapy of prelade to p



Logistic loss

Interpretation as binary cross-entropy loss

By interpreting the logistic function as probability of a label,

$$Z' = \omega_s + \omega_1 x_+ \omega_2 \times_2 + \cdots + \omega_c | \times_c |$$

$$y' = \frac{1}{1 + e^{2}};$$

The cross-entropy between the true (unknown) distribution and the estimated distribution

$$-\frac{1}{N}\left(\sum_{i=1}^{N}y^{i}\right)\log\frac{1}{1+e^{2i}}+\left(\left(-y^{i}\right)\log\frac{1}{1+e^{2i}}\right)$$

The cross entropy loss below, is the same as the logistic loss

=
$$\frac{1}{N}$$
 $\underset{i=1}{\overset{N}{\leq}}$ $\underset{i=1}{\overset{N}{\overset{N}{\leq}}$ $\underset{i=1}{\overset{N}{\leq}}$ $\underset{i=1}{\overset{N}{\leq}}$ $\underset{i=1}{\overset{N}{\simeq}}$ $\underset{i=1}{\overset{N}{$

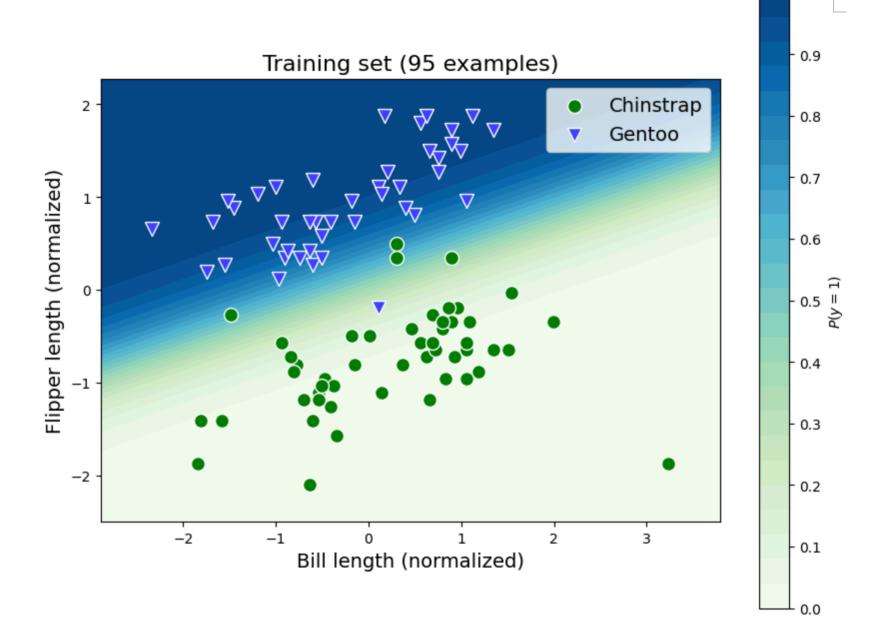


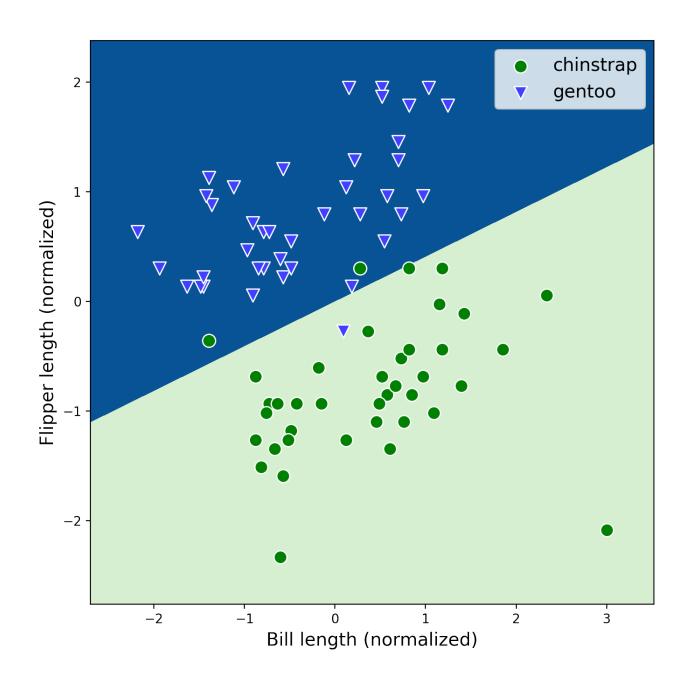
Logistic regression output

Palmer Penguins

	species	bill_length_mm	bill_depth_mm	flipper_length_mm	body_mass_g
0	Chinstrap	49.0	19.5	210.0	3950.0
1	Chinstrap	50.9	19.1	196.0	3550.0
2	Gentoo	42.7	13.7	208.0	3950.0
3	Chinstrap	43.5	18.1	202.0	3400.0
4	Chinstrap	49.8	17.3	198.0	3675.0



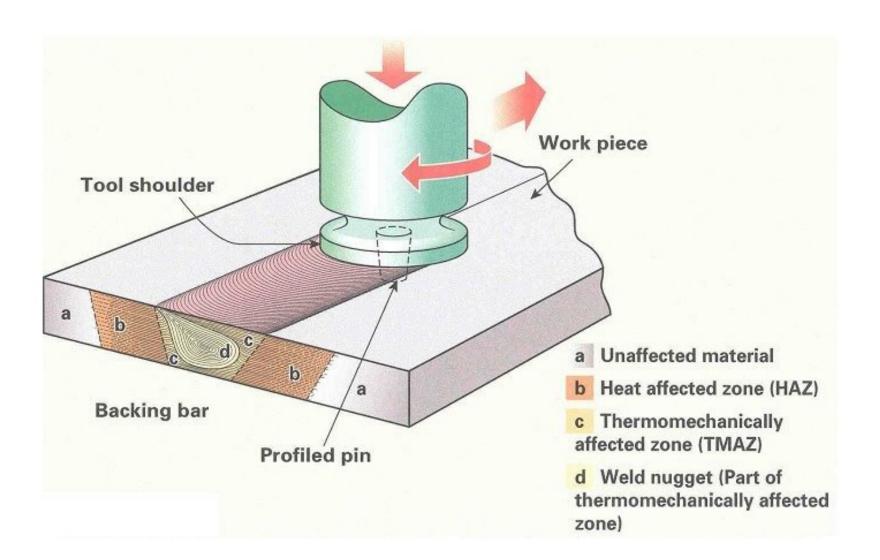


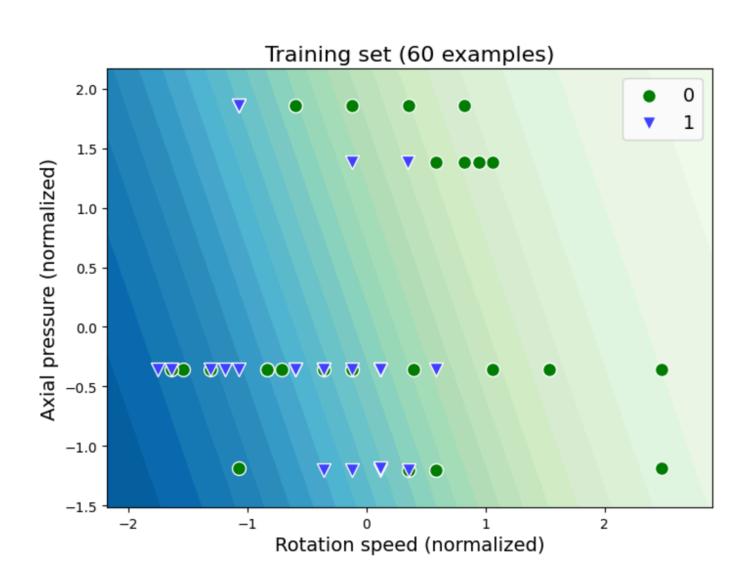




Logistic regression exercise this week

- Dataset 1: Void Formation in Welding, based on the paper
- Goal: formation of voids in friction stir welding as a function of the operation conditions
 - Tool rotational speed, axial pressure
 - The label: void or not void





- Dataset 2: discriminate between sonar signals bounced off a mine (metal cylinder) and those bounced off a roughly cylindrical rock
- Goal: predict whether the object is mine or rock based on
- The features (60 of them) are the energy within a particular frequency band, integrated over a certain period of time
- The label: rock/mine



Note on implementation Data normalization

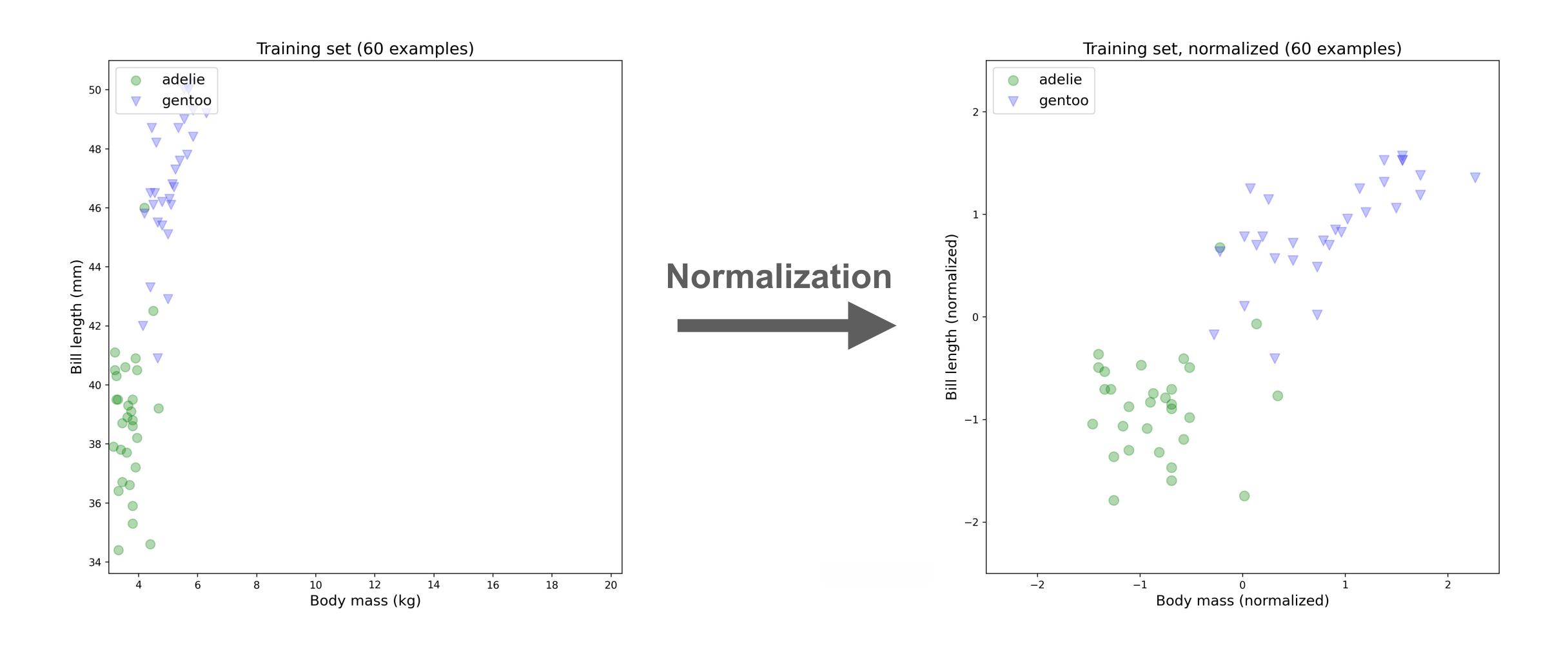
Data normalization / feature scaling: Normalize features (bring them all to the same scale)

Crucial step in preprocessing:

- Many classifiers (e.g. KNN that we will see in next lecture) rely on distance metrics
- Gradient descent will converge faster
- Coefficients are penalized appropriately (in the case where regularization is applied)



Numerical features Data normalization - Example





Performance metrics for binary classification Confusion matrix

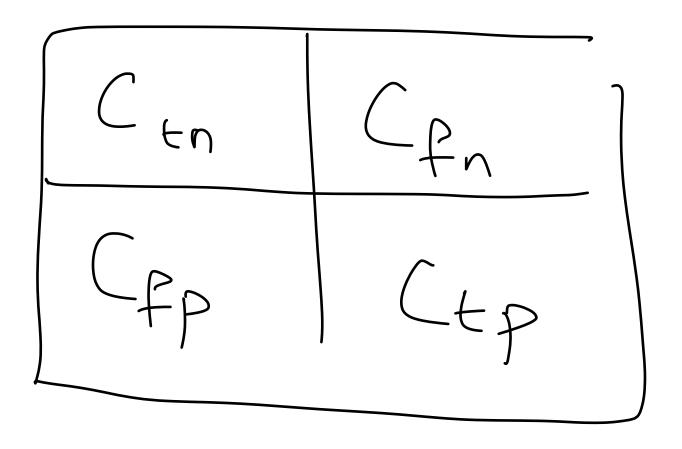
false positives: true label is 0, we predicted 7

false negatives: in 1, 1, 11 11 0

true positives: 11 11 11 11 11 11

true negatives: 11 11 11 0, 11 11 0

Confusion





Performance metrics for binary classification

Accuracy, error rate, recall

accuracy =
$$\frac{C_{tn} + C_{tp}}{N}$$

error rate = $\frac{C_{tn} + C_{tp}}{C_{tn} + C_{tp}}$



Exercise

Performance metric for binary classification

- We have used two approaches to train classifiers for spam email detection: "non-spam" (class 0) and "spam" (class 1)
- Our test set has 1000 emails, 900 of which were non-spam
- Approach 1: classified all data as non-spam
- Approach 2: classified 850 of non-spam emails as non-spam and 50 of spam emails as spam
- Write the confusion matrix of each approach
- Compute the error rate, accuracy and recall of each algorithm
- What do you conclude?



Outline

- Overfitting/underfitting and regularization
 - Connections to sensitivity and Lipschitz constant of the predictor
- Logistic regression
 - differentiable and convex loss function
 - Probabilistic interpretation
 - Performance metric different than the loss function
- Announcements:
- Exercise hour: logistic regression (notice data normalization more about it next week)
- Check moodle for past year's exams and quizzes