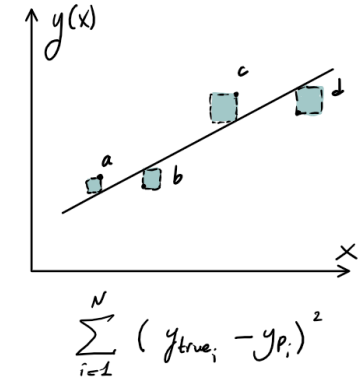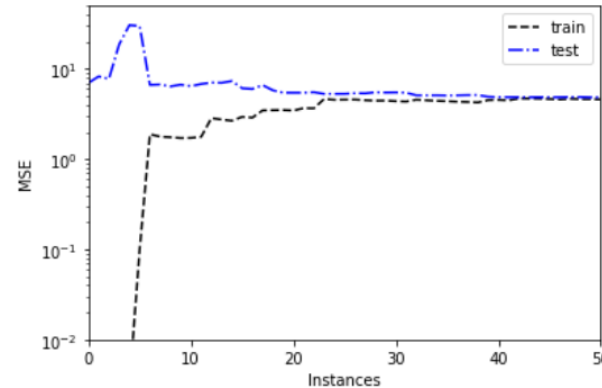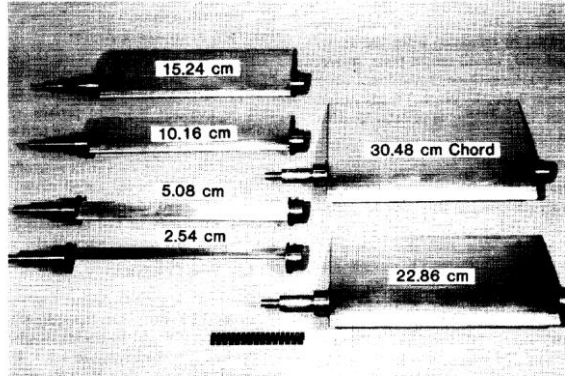# Data Driven Engineering I: Machine Learning for Dynamical Systems

**Analysis of Static Datasets I: Regression**

Institute of Thermal Turbomachinery
Prof. Dr.-Ing. Hans-Jörg Bauer

**www.its.kit.edu**

One page summary: Intro. to ML

* There are **4** main learning strategies, mainly based on *feedback* info.
  • Error-based, similarity-based, info-based, probab.-based

* The goals of **4** main ML tasks is very relevant to our learning strategies

* ML := ill-posed problem >> There will be many solutions for a problem.

* Nature & quality of data affects outcomes drastically !

* ML is very similar to cooking: Follow the proposed steps for a generic project.
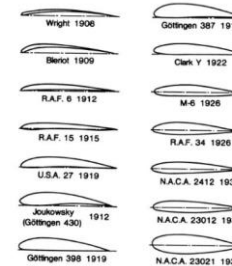
# Today's Agenda

Basic Steps to Follow:

0.) Understand the business/task.

1.) Understand the data.

2.) Explore & prepare the data.

3.) Shortlist candidate models.

4.) Training the model

5.) Evaluate the model predictions.

6.) "Serve" the model

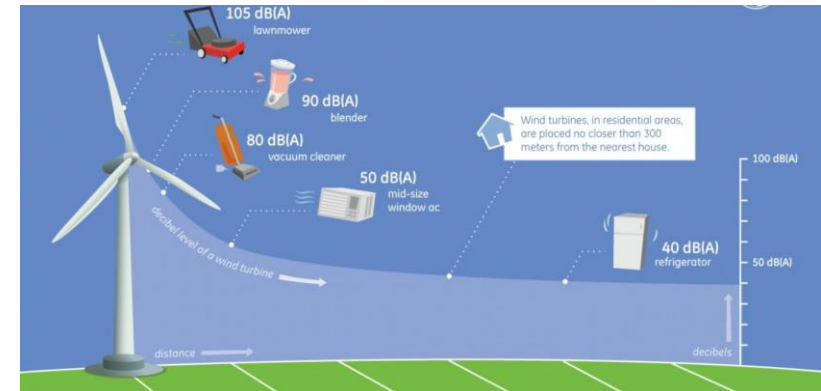Regression

# #0 Understanding the task

□ **Problem**: NACA 0012 Airfoil Noise Prediction based on Wind Tunnel Testing

□ **Noise** generated by an aircraft is an **economic** (efficiency) and **enviromental** issue.

□ One component of the noise the **self-noise of the airfoil**: interaction of the airfoil with its own boundary layer

1917, the NACA Technical Report No. 18 titled "Aerofoils and Aerofoil Structural Combinations," was released.

# #0 Understanding the task

❑ Engineering: semi-emprical models (Brooks)

❑ Five self-noise mechanisms due to specific boundary-layer phenomena have been identified

❑ The database is from seven NACA0012 airfoil blade sections of different sizes tested at wind tunnel speeds up to Mach 0.21 and at angles of attack from 0°to 25.2°.

  ✓ Freq. of noise
  ✓ Angle of attack
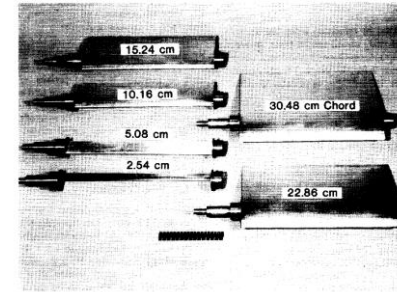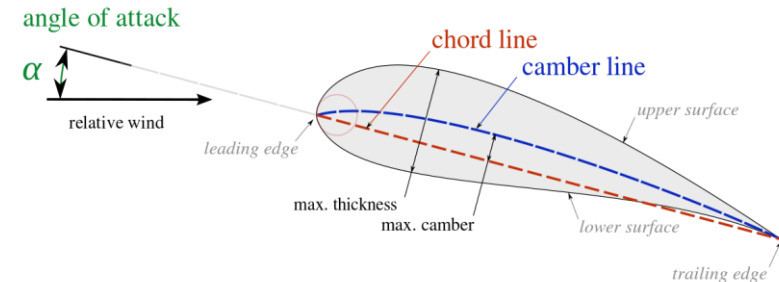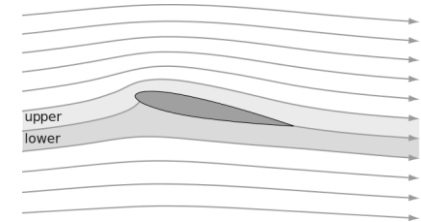  ✓ Free stream velocity
  ✓ Geometry of the airfoil



Figure 2. Two-dimensional NACA 0012 airfoil blade models.





angle of attack

$\alpha$

relative wind

leading edge

chord line

camber line

upper surface

max. thickness

max. camber

lower surface

trailing edge

# #1 Understanding the data

☐ Check the data source: understand what the data refers to

☐ Objective: understand the characteristics of the data

☐ Look at the feature columns:
  - ☐ Any missing values?
  - ☐ Any features with NaN values?
  - ☐ Uniqueness of the dataset? ("cardinality")

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1503 entries, 0 to 1502
Data columns (total 6 columns):
 #   Column                 Non-Null Count   Dtype
---  ------                 --------------   -----
 0   frequency              1503 non-null    int64
 1   angle_attack           1503 non-null    float64
 2   chord_length           1503 non-null    float64
 3   Free-stream_velocity   1503 non-null    float64
 4   displacement_thickness 1503 non-null    float64
 5   sound_pressure         1503 non-null    float64
dtypes: float64(5), int64(1)
memory usage: 70.6 KB
```

```
data.head(5)
```

| | frequency | angle_attack | chord_length | Free-stream_velocity | displacement_thickness |
|---|---|---|---|---|---|
| 0 | 800 | 0.0 | 0.3048 | 71.3 | 0.002663 |
| 1 | 1000 | 0.0 | 0.3048 | 71.3 | 0.002663 |
| 2 | 1250 | 0.0 | 0.3048 | 71.3 | 0.002663 |
| 3 | 1600 | 0.0 | 0.3048 | 71.3 | 0.002663 |
| 4 | 2000 | 0.0 | 0.3048 | 71.3 | 0.002663 |

Ateliers & Saveurs in Montreal
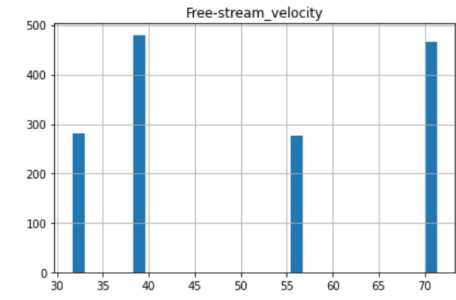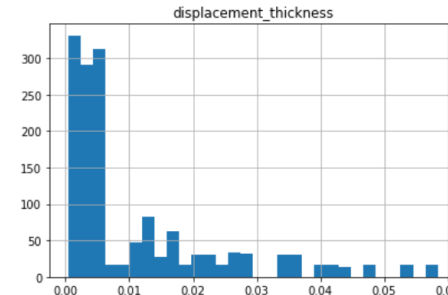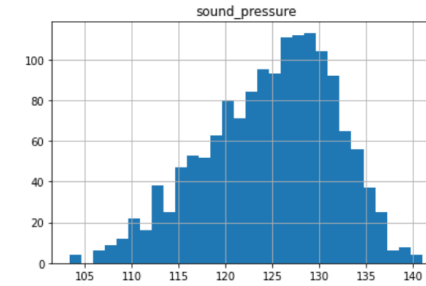
Institute of Thermal Turbomachinery (ITS)

# #2 Exploring the data

❑ **Objective**: generate a data quality report

❑ Using standard statistical measures of central tendency and variation
  ❑ tabular data and visual plots
  ❑ mean, mode, and median
  ❑ standard deviation and percentiles
  ❑ bars, histograms, box and violin plots

✓ Missing values,
✓ Irregular cardinality problems,
  ▪ 1 or comparably small
✓ Outliers
  ▪ invalid outliers and valid outliers

# #2 Exporing the data: Correlation Matrix

❑ Shows the correlation between each pair of features

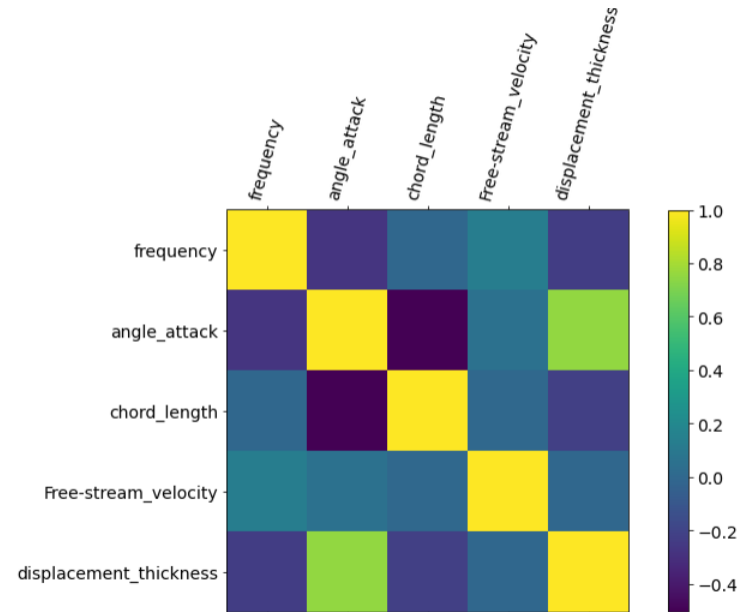$$Cov(a,b) = \frac{1}{n-1} \sum_{i=1}^{n} \left[ (a_i - \bar{a}) \times (b_i - \bar{b}) \right]$$

$\downarrow \downarrow$ Features

$\downarrow$ instance

$\downarrow$ mean

$\downarrow$ mean

❑ Normalized form of "covariance"

$$Corr(a,b) = \frac{Cov(a,b)}{SD(a) \times SD(b)}$$

}
* Normalized
* Dimensionless
Easy to interpret

❑ Ranges between −1 and +1

Ateliers & Saveurs in Montreal

Institute of Thermal Turbomachinery (ITS)

# #2 Preparing the Data

❑ Classification >> supervised >> **training & test split**



❑ Reducing overfitting via **cross-validation**: take **random portions** of the data to build a model

❑ **k-fold** method: k = 5; (typically 10)

Ateliers & Saveurs in Montreal

colab

# # Model Selection: Linear Regression 1

```
data.head(5)
```

| | frequency | angle_attack | chord_length | Free-stream_velocity | displacement_thickness |
|---|---|---|---|---|---|
| 0 | 800 | 0.0 | 0.3048 | 71.3 | 0.002663 |
| 1 | 1000 | 0.0 | 0.3048 | 71.3 | 0.002663 |
| 2 | 1250 | 0.0 | 0.3048 | 71.3 | 0.002663 |
| 3 | 1600 | 0.0 | 0.3048 | 71.3 | 0.002663 |
| 4 | 2000 | 0.0 | 0.3048 | 71.3 | 0.002663 |

$$X \rightarrow \boxed{\quad} \rightarrow y$$

Model

$$y \approx f(X, w)$$

features $X$

labels $y$

What we know     What we want

①  $y_p = w_0 + w_1 X_1 + w_2 X_2 + \dots w_n X_n$  } linear func. $w_i$ to $x_i$

$\downarrow$ (bias)

②  $y_{True_i} = y_{P_i} + Error_i$  } Error Metric (norm) $:\approx$ Goodness of a fit

③  Extended to Nonlinearity via $\phi$

$$y_p = w_0 + \sum w_i \phi(x_i)$$ $\rightarrow$ Basis functions $\rightarrow x^i$ (polynomial)

linear   nonlinear     $\dots$ $\rightarrow \sigma\left(\frac{x-\mu}{s}\right)$ (sigmoidal)

# # Model Selection: Linear Regression 2



$x \rightarrow$ [Model] $\rightarrow y$

$\textcircled{2}$  $y_{True_i} = y_{p_i} + Error_i$ $\Big\}$  Error Metric (norm) $:\simeq$ Goodness of a fit

- Maximum error  $(\ell_\infty)$  $\max\limits_{1 \leq i \leq n} \left| y_{true_i} - y_{p_i} \right|$

- Mean absolute error  $(\ell_1)$  $\frac{1}{n} \sum\limits_{i=1}^{n} \left| y_{true_i} - y_{p_i} \right|$

- Least Squares error  $(\ell_2)$  $\left( \frac{1}{n} \sum\limits_{i=1}^{n} \left| y_{true_i} - y_{p_i} \right|^2 \right)^{1/2}$

Error-based learning:

$y_{True_i} = y_{p_i} + Error_i$

$X_2$

$X_1$

$e_n$

$e_1$
$e_2$

outlier ⚠

Presence of outliers / limited observ.

$$\max_{1 < i < n} | y_{true_i} - y_{p_i} |$$

$$\frac{1}{n} \sum_{i=1}^{n} | y_{true_i} - y_{p_i} |$$

$$\left( \frac{1}{n} \sum_{i=1}^{n} | y_{true_i} - y_{p_i} |^2 \right)^{1/2}$$

Errors will be dictated by outliers ⚠

💡 Regularization

# Model Selection: Linear Regression 3



④ Regularized Linear Regression

O'is // "Over-fitted" Regularized $\Rightarrow$ $E_T = (E_D) + (E_R)$ Regularization Error

Data error

$E_R \leftarrow \dfrac{\lambda}{2} \sum\limits_{i}^{M} |w_i|^q$

ⓐ Ridge Regression $\Rightarrow$ $E_R = \dfrac{\alpha}{2} \sum\limits_{i=1}^{n} w_i^2$ (no bias here)

ⓑ Lasso $\Rightarrow$ $E_R \Leftrightarrow \ell_1$ norm $E_R = \dfrac{\alpha}{2} \sum\limits_{i=1}^{n} |w_i|$ $\left( \begin{array}{l} \alpha \text{ is large;} \\ w \text{ is sparse} \end{array} \right)$

(1-r) Ridge     (r) Lasso

ⓒ Elastic Net $\Rightarrow$ $E_R = \underbrace{\dfrac{1-r}{2} \alpha \sum w_i^2}_{} + \underbrace{\dfrac{r\alpha}{2} \sum |w_i|}_{}$

(i) Large $\lambda$

"Lasso"

$E_R \propto (w_1 + w_2)$

# #4 Training the model

☐ Classification >> supervised >> **training & test split**



☐ Reducing overfitting via **cross-validation**: take **random portions** of the data to build a model

☐ **k-fold** method: k = 5; (typically 10)



Institute of Thermal Turbomachinery (ITS)

# #5 Evaluation of the results

❑ Coefficient of determination, $R^2$

- Indicates the goodness of fit
- Measure of generalization capability
- Best possible score is 1.0
- It can be negative



$$\sum_{i=1}^{N} \left( y_{true_i} - y_{p_i} \right)^2$$

$$\sum_{i=1}^{N} \left( y_{true_i} - \bar{y} \right)^2$$

$$R^2 \left( y_{true}, y_p \right) = 1 - \frac{\sum_{i=1}^{N} \left( y_{true_i} - y_{p_i} \right)^2}{\sum_{i=1}^{N} \left( y_{true_i} - \bar{y} \right)^2} \qquad \bar{y} = \frac{1}{N} \sum_{i=1}^{N} y_{true_i}$$

Ateliers & Saveurs in Montreal

Case: 'Curve Fitting'

Linear
Polynomial (n=2)
Polynomial (n=7)

(i) model learns
(ii) as it learns, model parameters generalizes.
(iii) $E_D$ is found.

(i) Compare it with n=7;
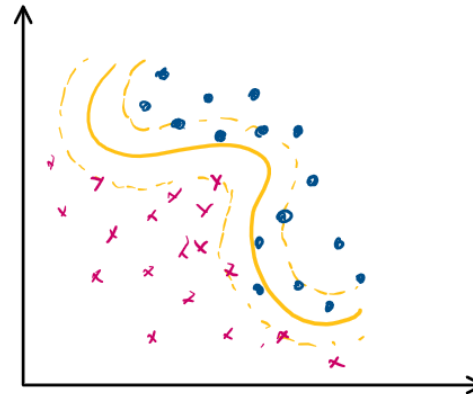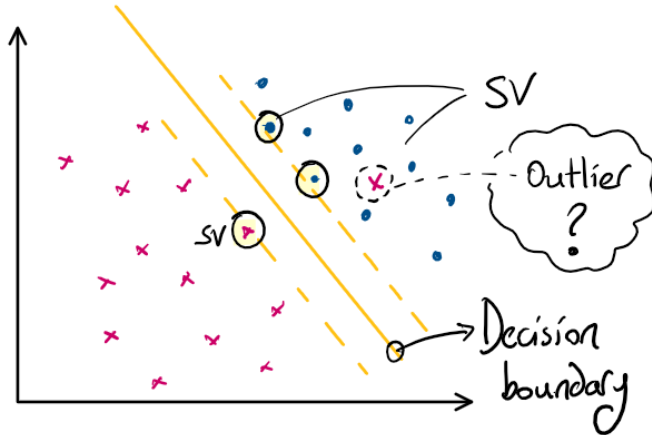(ii) Divergence of $E_D$ ⟹ Overfitting

Ateliers & Saveurs in Montreal

— SVM —

Classification

* fits a "street" between classes
* uses support vectors (sv)
* Decision is based on SVs, not other instances.
* Feature scaling is important
* outliers ⟹ "Soft Margin" (~c)
  ✓ limit margin violations
* must be linearly separable

SV
Outlier?
Decision boundary

— SVM —

Classification



SV

Outlier ?

SV

Decision boundary

$*$ linear decision bound. $\Rightarrow$ X

$*$ "Kernel Trick" := $\phi(x)$

($\checkmark$) introduce non-linearity

($\checkmark$) "feature eng." without adding new features.

— SVM —

Regression

* Fit as many instance as possible

* "Street" width is controlled by margin $\epsilon$.

* Convex optimization problem:

☑ C

☑ $\epsilon$

☑ Kernel

$$- SVM -$$

☑ C
☑ $\epsilon$
☑ Kernel

① $E_T = \boxed{E_D} + E_R$

$\longrightarrow$ Replaced by an $\epsilon$-insensitive function:

② $E_D = \begin{cases} 0, & \text{if } |y_{true} - y_p| < \epsilon \\ |y_{true} - y_p| - \epsilon, & \text{otherwise} \end{cases}$

③ We minimize:  Kernel $\quad \epsilon$

$\boxed{C} \sum_{i=1}^{N} \left[ |y_{true_i} - y_{p_i}| - \epsilon \right] + \frac{1}{2} \frac{1}{N} \sum_{i=1}^{N} w_i^2$

$\downarrow$ Regularization parameter

Ateliers & Saveurs in Montreal

# Model Selection: Bayesian Regression 1

# Model Selection: Bayesian Regression 2

① Bayesian approach; $y_t = y_p + \underline{Error}$

                                    "Gaussian noise"

② $p(y_t | X, w, \alpha) = \mathcal{N}(y_t | y_p, \alpha) \Rightarrow \alpha$

            "Given that"

                                     Hyperparameters

③ $p(w | \lambda) = \mathcal{N}(w | 0, \lambda^{-1} I_p) \Rightarrow \lambda$     in scikit learn ⚡

Bayes' Theorem

* "Hypothesis" + "evidence" = "New hypothesis"

# Model Selection: Bayesian Regression 3

Bayes' Theorem

\* "Hypothesis" + "evidence" = "New hypothesis"
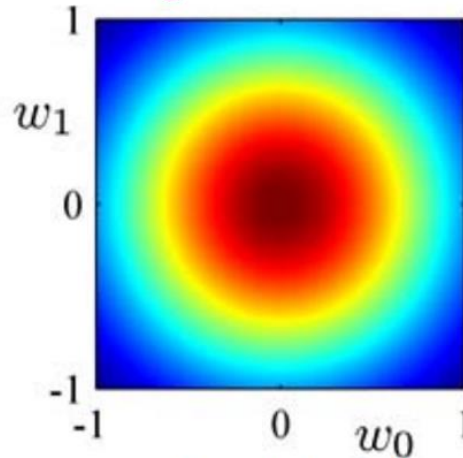
Prior knowledge $\Longrightarrow$ Posterior knowledge
$\Downarrow$
"probability"

\* $\text{Posterior prob.} = \dfrac{\text{Likelihood} \times \text{Prior Knowledge or prob.}}{\text{Evidence}}$
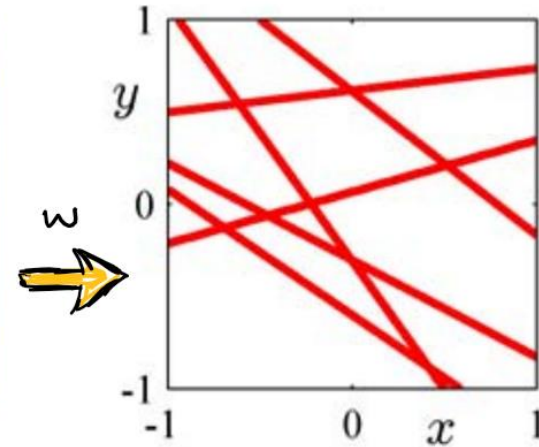
Before any data points observed:

Model:

$$y = w_0 + w_1 x$$

prior distribution of $w$

$y(x, w)$

Pattern Recognition and Machine Learning, Chapter 3

After observing 1 data point:

True value

$w_1$

$w_0$

$w_1$

$w_0$

$w$

$y$

Data

$x$

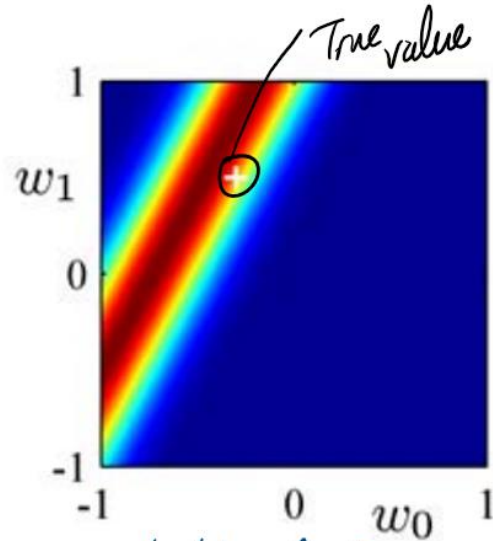Likelihood function:

* $P(y_t / x, w)$

* $w \leftarrow$ "Prior"

Posterior (updated) probability
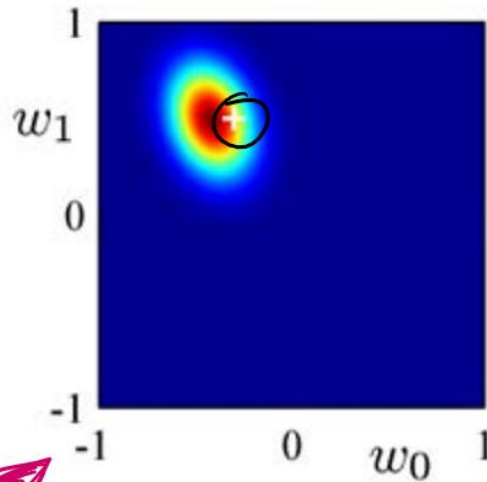
$y(x, w)$

* Lines are being accumulated around data.

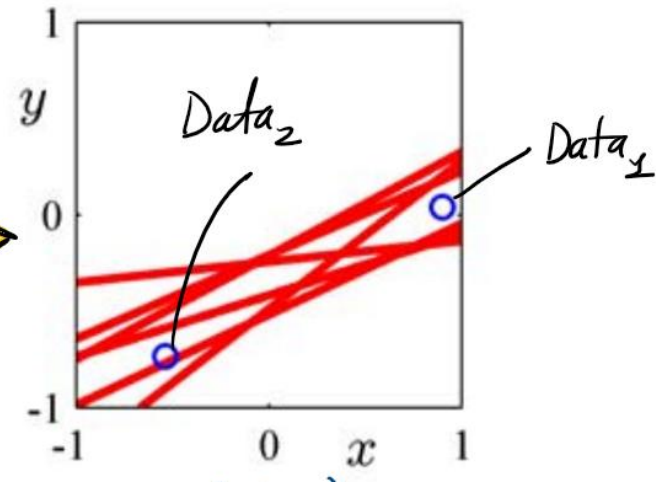Pattern Recognition and Machine Learning, Chapter 3

Effect of Observing the second data:

True value

$w_1$

$w_0$

Likelihood func.:
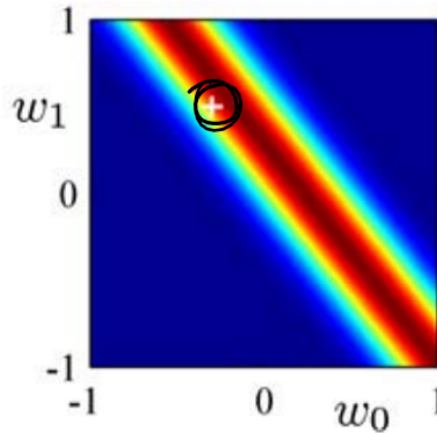$p(y_t \mid x, w)$

$w_1$

$w_0$

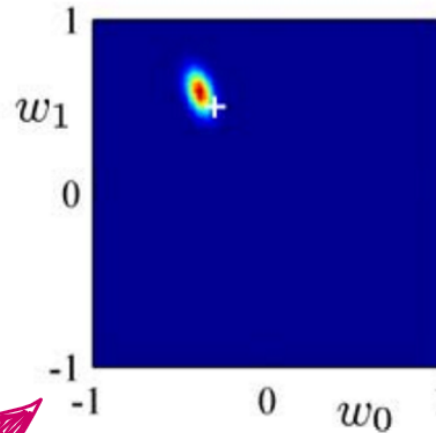Posterior probability
(updated)

$w$

$y$

$Data_2$

$Data_1$

$x$

$y(x, w)$

* lines are accum.

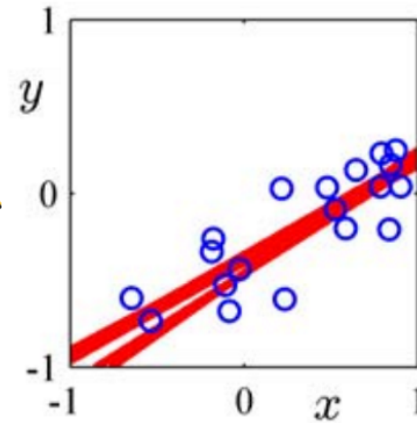Pattern Recognition and Machine Learning, Chapter 3

Effect of observing 20 data points

Likelihood func.

Posterior probability (updated)

$y(x, w)$

* Lines are condensed.

Pattern Recognition and Machine Learning, Chapter 3

Ateliers & Saveurs in Montreal

# Additional Notes