

Lecture 7a

**Prediction**

# Outline

## Motivations

### Least Square Regression

- ▶ Model and Objective

### Connection to CCA

- ▶ Least square regression as a special case of CCA

### Computational Aspects

- ▶ Problem of invertibility
- ▶ High-dimensions

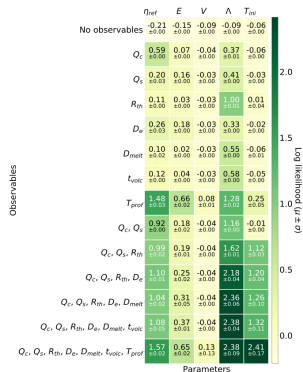
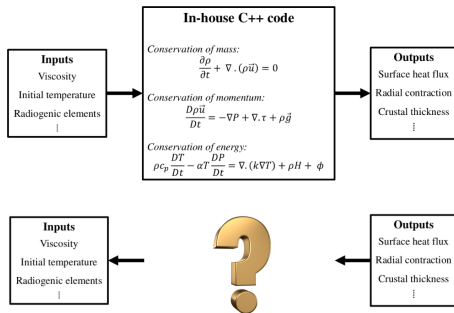
### Support Vector Regression

- ▶ Beyond least square regression
- ▶ Absolute deviations and  $\epsilon$ -insensitivity
- ▶ SVR primal
- ▶ SVR dual

# Motivations

## Example: Space Science

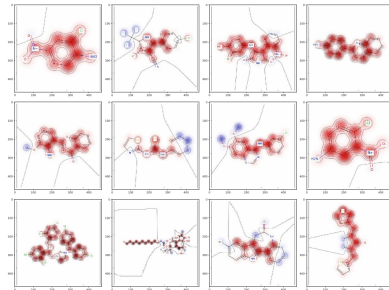
- To what extent can physical parameters of (exo-)planets (e.g. viscosity) be inferred from observables [1].



# Motivations

## Example: Toxicity prediction

- Are particular substructures predictive of molecular toxicity, and can we predict toxicity from the molecular geometry [3].



Image

[https://doi.org/10.1007/978-3-030-28954-6\\_18](https://doi.org/10.1007/978-3-030-28954-6_18)

source:

Part 1

# Least Squares Regression

# Least Squares Regression

Idea:

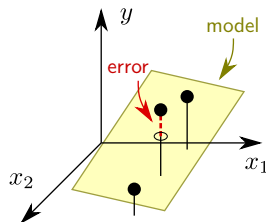
- ▶ Build a model predicting  $y$  from  $\mathbf{x}$ :

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$$

- ▶ Find model parameters so that the prediction errors are minimized, e.g. least square error averaged over all instances:

$$\min_{\mathbf{w}, b} \mathbb{E}[(\mathbf{w}^\top \mathbf{x} + b - y)^2]$$

- ▶ The error of the model can either be of interest by itself (quantifying predictability), or the learned parameters can be inspected for further insights.



# Simplification for Centered Data

## Proposition

*Assume we wish to train a linear model*

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$$

*on data that has been centered as a first step, i.e.  $E[\mathbf{x}] = \mathbf{0}$  and  $E[y] = 0$ . Then, we can show that*

$$\arg \max_b E[(f(\mathbf{x}) - y)^2] = 0$$

*i.e. it is always best to set the bias to zero.*

**Proof:**

$$\frac{\partial}{\partial b} E[(f(\mathbf{x}) - y)^2] = 2E[(\mathbf{w}^\top \mathbf{x} + b - y)] \stackrel{\text{def}}{=} 0 \Rightarrow b = 0$$



**Implication:**

- ▶ Without loss of accuracy, one can simplify the least square regression problem as finding a homogeneous linear model  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$  that minimizes the mean square error.

# Linear Regression

- Recall that we consider prediction functions of the type  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$  and we would like to minimize the mean square error. The latter can be developed as:

$$\begin{aligned}\mathcal{E}(\mathbf{w}) &= \mathbb{E}[(f(\mathbf{x}) - y)^2] \\ &= \mathbb{E}[(\mathbf{w}^\top \mathbf{x} - y)^2] \\ &= \mathbb{E}[\mathbf{w}^\top \mathbf{x} \mathbf{x}^\top \mathbf{w} - 2\mathbf{w}^\top \mathbf{x} y + y^2] \\ &= \mathbf{w}^\top C_{xx} \mathbf{w} - 2\mathbf{w}^\top C_{xy} + C_{yy}\end{aligned}$$

- Observing that  $C_{xx}$  is positive semi-definite, minimizing  $\mathcal{E}(\mathbf{w})$  is a convex optimization problem and the solution must necessarily have gradient zero, i.e.  $\nabla \mathcal{E}(\mathbf{w}) = 2C_{xx} \mathbf{w} - 2C_{xy} = 0$ . Solving for  $\mathbf{w}$  gives us the optimal model:

$$\boxed{\mathbf{w} = C_{xx}^{-1} C_{xy}}$$

- Injecting this solution into the objective, we get the mean square error at the optimum:

$$\mathcal{E} = C_{yy} - C_{yx} C_{xx}^{-1} C_{xy}$$



# Example of a Linear Model

Model	Unstandardized Coefficients		Standardized Coefficients	t	Sig.	95.0% Confidence Interval for B	
	① B	Std. Error	③ Beta			Lower Bound	Upper Bound
1 (Constant)	-3263.586	1059.163		-3.081	.002	-5344.360	-1182.812
Sex	509.271	139.565	.126	3.649	.000	235.089	783.454
Age at Survey Completion (Years)	114.658	12.452	.322	9.208	.000	90.196	139.121
Average Consumption of Alcoholic Beverages per Week	50.386	10.275	.192	4.904	.000	30.201	70.572
Average Consumption of Cigarettes per Day	139.414	17.384	.311	8.020	.000	105.263	173.565
Average Hours of Exercise per Week	-271.270	36.300	-.281	-7.473	.000	-342.584	-199.956

a. Dependent Variable: Total Health Care Costs Declared over 2020

www.spss-tutorials.com

Image source: <https://www.spss-tutorials.com/spss-multiple-linear-regression-example/>

## Remarks:

- ▶ Model weights can be interpreted as the amount by which the output would increase if changing the input variable by +1. To assess the relative importance of variables, it is necessary to standardize the input features.
- ▶ Unless the input data are decorrelated ( $C_{xx}$  diagonal), it is not impossible that input features that positively correlate with the output contribute *negatively* in the linear model. → one must be careful with interpreting the weights of a linear model.

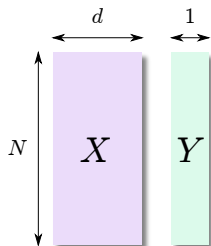
Part 2

## **Connection to CCA**

# Connection to CCA

- ▶ Recall that canonical correlation analysis aims to find a projection of two modalities that maximizes their correlation, and that such correlation can be expressed in terms of auto- and cross-covariances as:

$$\rho = \frac{\mathbf{w}_x^\top C_{xy} \mathbf{w}_y}{\sqrt{\mathbf{w}_x^\top C_{xx} \mathbf{w}_x} \sqrt{\mathbf{w}_y^\top C_{yy} \mathbf{w}_y}}$$



- ▶ Because in the case of linear regression the second modality is univariate, there is no direction to be found in the second modality (we can set manually  $w_y = 1$ ), and the CCA objective simplifies to:

$$\rho = \frac{\mathbf{w}_x^\top C_{xy}}{\sqrt{\mathbf{w}_x^\top C_{xx} \mathbf{w}_x} C_{yy}}$$

- ▶ This can be reformulated as a constrained optimization problem:

$$\max_{\mathbf{w}} \mathbf{w}_x^\top C_{xy} \quad \text{s.t.} \quad \mathbf{w}_x^\top C_{xx} \mathbf{w}_x C_{yy} = 1$$

# Connection to CCA: Deriving the Weights

- ▶ We would to solve the constrained optimization problem:

$$\max_{\mathbf{w}} \mathbf{w}_x^\top C_{xy} \quad \text{s.t.} \quad \mathbf{w}_x^\top C_{xx} \mathbf{w}_x C_{yy} = 1$$

- ▶ Using the framework of Lagrange multipliers, we get:

$$\mathcal{L}(\mathbf{w}_x, \lambda) = \mathbf{w}_x^\top C_{xy} + \frac{1}{2} \lambda (1 - \mathbf{w}_x^\top C_{xx} \mathbf{w}_x C_{yy})$$

and

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}_x} = C_{xy} - \lambda C_{xx} \mathbf{w}_x C_{yy} \stackrel{!}{=} 0 \quad \Rightarrow \quad \mathbf{w}_x = \lambda^{-1} C_{xx}^{-1} C_{xy} C_{yy}^{-1}$$

and setting  $\lambda$  in a way that the constraint is satisfied, we get:

$$\mathbf{w}_x = \frac{C_{xx}^{-1} C_{xy}}{\sqrt{C_{yx} C_{xx}^{-1} C_{xy} C_{yy}}}.$$

- ▶ This is exactly the same direction as the weight of the linear model ( $\mathbf{w} = C_{xx}^{-1} C_{xy}$ ).

# Connection to CCA: Deriving the Objective

- ▶ Recall that the parameter that optimizes our CCA objective is given by:

$$\mathbf{w}_x = \frac{C_{xx}^{-1} C_{xy}}{\sqrt{C_{yx} C_{xx}^{-1} C_{xy} C_{yy}}}$$

- ▶ Evaluating the CCA objective with this solution gives the correlation coefficient

$$\rho = \mathbf{w}_x^\top C_{xy} = \sqrt{C_{yx} C_{xx}^{-1} C_{xy} C_{yy}^{-1}}.$$

From this correlation coefficient, one can measure the explained variance  $C_{yy} \rho^2 = C_{yx} C_{xx}^{-1} C_{xy}$ .

- ▶ Compare with the error of our linear model that was given by:

$$\mathcal{E} = C_{yy} - C_{yx} C_{xx}^{-1} C_{xy}$$

- ▶ The objectives of CCA and least square regression are therefore related as  $C_{yy} = C_{yy} \rho^2 + \mathcal{E}$ , i.e. together they form a decomposition of the total variance (i.e. maximizing the correlation using CCA or minimizing the mean square error achieve both the same objective!).

Part 3

## **Computational Aspects**

# Back to Linear Regression ...

Recall:

- ▶ When considering prediction functions of the type  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$ , and minimizing the prediction error

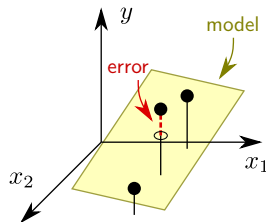
$$\mathcal{E}(\mathbf{w}) = \mathbb{E}[(f(\mathbf{x}) - y)^2].$$

we get the optimal model

$$\mathbf{w} = C_{xx}^{-1} C_{xy}$$

and its error is:

$$\mathcal{E} = C_{yy} - C_{yx} C_{xx}^{-1} C_{xy}$$



Question:

- ▶ What if  $C_{xx}$  is not invertible?

# Dealing with Non-Invertible $C_{xx}$

## Observation:

- ▶  $C_{xx}$  non-invertible corresponds to the case where the data only spans a subspace of  $\mathbb{R}^d$ . In that case, there are infinitely many possible solutions to the regression problem.

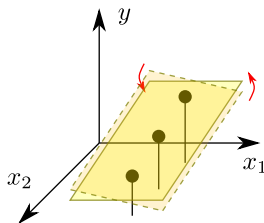
## Idea:

- ▶ (Virtually) add some small uncorrelated noise  $\mathbf{n}$  to the data.

$$\begin{aligned}C_{xx}^{\text{new}} &= \text{Cov}(\mathbf{x} + \mathbf{n}, \mathbf{x} + \mathbf{n}) \\&= \text{Cov}(\mathbf{x}, \mathbf{x}) + 2\text{Cov}(\mathbf{x}, \mathbf{n}) + \text{Cov}(\mathbf{n}, \mathbf{n}) \\&= C_{xx} + \sigma_n^2 I\end{aligned}$$

This corresponds to adding a diagonal term to the covariance matrix.

- ▶ This makes the covariance matrix invertible and favors solutions that are flat on directions that are orthogonal to the data.





# Dealing with Non-Invertible $C_{xx}$

## Implication on the model error:

- ▶ Let us virtually inject the uncorrelated noise in the data, obtain the resulting robustified covariance matrix  $C_{xx}^{\text{new}} = C_{xx} + \sigma_n^2 I$ , and also observe that the covariance matrices  $C_{xy}, C_{yy}$  are not affected by that noise.
- ▶ Using the new covariances in the error function yields:

$$\begin{aligned}\mathcal{E}^{\text{new}} &= C_{yy} - C_{yx} (C_{xx}^{\text{new}})^{-1} C_{xy} \\ &= C_{yy} - C_{yx} (C_{xx} + \sigma_n^2 I)^{-1} C_{xy} \\ &= C_{yy} - C_{yx} (\sum_{j=1}^d u_j u_j^\top (\lambda_j + \sigma_n^2))^{-1} C_{xy} \\ &= C_{yy} - \sum_j (C_{yx} u_j)^2 \frac{1}{\lambda_j + \sigma_n^2}\end{aligned}$$

The larger the noise  $\sigma_n^2$ , the higher the error (and the lower the correlation).  
→ Just use  $\sigma_n$  large enough to be able to invert  $C_{xx}$ .

## Note:

- ▶ This is not the whole picture yet, a higher value of  $\sigma_n^2$  can be beneficial for *reproducibility* of results (→ next week).

# Dealing with High Dimensions

- ▶ When  $d$  is large, computing and inverting the matrix  $C_{xx}$  can be prohibitively expensive.
- ▶ Let us start from the original formulation of the regression objective

$$\begin{aligned}\mathcal{E}(\mathbf{w}) &= \mathbb{E}[(\mathbf{w}^\top \mathbf{x} - y)^2] \\ &= \mathbf{w}^\top C_{xx} \mathbf{w} - 2\mathbf{w}^\top C_{xy} + C_{yy}\end{aligned}$$

- ▶ We now assume that an optimal  $\mathbf{w}$  is found in the span of the data, and express  $\mathbf{w}$  as  $\mathbf{w} = X\boldsymbol{\alpha}$ , a linear combination of the data, where  $X$  is a matrix of size  $d \times N$  containing the centered data:

$$\begin{aligned}\mathcal{E}(\boldsymbol{\alpha}) &= \overbrace{\boldsymbol{\alpha}^\top X^\top}^{\mathbf{w}^\top} C_{xx} \overbrace{X\boldsymbol{\alpha}}^{\mathbf{w}} - 2\overbrace{\boldsymbol{\alpha}^\top X^\top}^{\mathbf{w}^\top} C_{xy} + C_{yy} \\ &= \frac{1}{N} \boldsymbol{\alpha}^\top \underbrace{X^\top X X^\top X}_{Q_x^2} \boldsymbol{\alpha} - \frac{2}{N} \boldsymbol{\alpha}^\top \underbrace{X^\top X Y}_{Q_x} + C_{yy}\end{aligned}$$

- ▶ No need to compute matrices of size  $d \times d$ . Matrices  $Q_x^2$  and  $Q_x$  are of size  $N \times N$ .

# Dealing with High Dimensions

- ▶ We have found that the least square error can be expressed as

$$\mathcal{E} = \frac{1}{N} \boldsymbol{\alpha}^\top Q_x^2 \boldsymbol{\alpha} - \frac{2}{N} \boldsymbol{\alpha}^\top Q_x Y + C_{yy}$$

- ▶ Observing that  $\mathcal{E}(\boldsymbol{\alpha})$  is convex, we find the solution where the gradient is zero:

$$\nabla \mathcal{E}(\boldsymbol{\alpha}) = \frac{2}{N} Q_x^2 \boldsymbol{\alpha} - \frac{2}{N} Q_x Y \stackrel{!}{=} 0$$

and we find the solution

$$\boldsymbol{\alpha} = (Q_x^2)^{-1} Q_x Y = Q_x^{-1} Y$$

Therefore, we only need to invert a matrix of size  $N \times N$ .

- ▶ From this solution, one can recover the original weight parameter as:

$$\boldsymbol{w} = X \boldsymbol{\alpha}$$

## Part 4

# Regression with Outliers

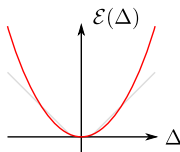
# Regression with Outliers

## Motivations:

- ▶ Square errors may be very large for data points whose outputs strongly deviate from that of other data points.
- ▶ Points with large prediction errors may be better treated as outliers (for which we accept the model to be inaccurate), so that the model can focus on the non-outlier part of the data.
- ▶ Reduced sensitivity to outliers can be achieved by considering absolute deviations instead of square errors, and invariance to small noise can be maintained by introducing a small slack  $\epsilon$ .

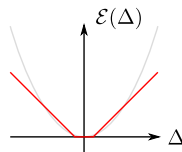
squared error

$$\mathcal{E}(\Delta) = \Delta^2$$



$\epsilon$ -insensitive absolute deviation

$$\mathcal{E}(\Delta) = \max(0, |\Delta| - \epsilon)$$



# Regression with Outliers

- ▶ Consider the original least square regression problem

$$\min_{\mathbf{w}} \mathbb{E}[\underbrace{(\mathbf{w}^\top x - y)^2}_{\Delta}]$$

- ▶ To reduce exposure to outliers, we can replace the square error by the ( $\epsilon$ -insensitive) absolute deviation.

$$\min_{\mathbf{w}, b} \mathbb{E}[\max(0, \underbrace{|\mathbf{w}^\top x + b - y|}_{\Delta} - \epsilon)]$$

(note that we need to reintroduce the bias, because we are no longer guaranteed that a solution without bias is optimal).

- ▶ There may be multiple solutions that solve the problem exactly (especially when  $d > N$ ), hence, we introduce a small penalty term  $\lambda \|\mathbf{w}\|^2$  that favors the flattest solution:

$$\min_{\mathbf{w}, b} \mathbb{E}[\max(0, \underbrace{|\mathbf{w}^\top x + b - y|}_{\Delta} - \epsilon)] + \lambda \|\mathbf{w}\|^2$$

# Support Vector Regression

An algorithm that implements these ideas is Support Vector Regression [2, 4]. It takes the form of a quadratic optimization problem with linear constraints.

Support Vector Regression (Primal):

$$\min_{\mathbf{w}, b, \xi, \xi^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N (\xi_i + \xi_i^*)$$

$$\begin{aligned} \text{s.t.} \quad & \forall_{i=1}^N : \mathbf{w}^\top \mathbf{x}_i + b - y_i \leq \epsilon + \xi_i \\ & \forall_{i=1}^N : y_i - \mathbf{w}^\top \mathbf{x}_i + b \leq \epsilon + \xi_i^* \\ & \xi_i, \xi_i^* \geq 0 \end{aligned}$$

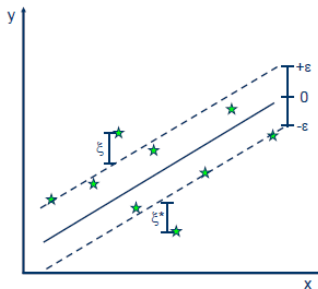


Image source: [https://www.saedsayad.com/support\\_vector\\_machine\\_reg.htm](https://www.saedsayad.com/support_vector_machine_reg.htm)

# SVR: Primal vs. Dual

Primal:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi, \xi^*} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N (\xi_i + \xi_i^*) \\ \text{s.t.} \quad & \forall_{i=1}^N : \mathbf{w}^\top \mathbf{x}_i + b - y_i \leq \epsilon + \xi_i \\ & \forall_{i=1}^N : y_i - \mathbf{w}^\top \mathbf{x}_i + b \leq \epsilon + \xi_i^* \\ & \xi_i, \xi_i^* \geq 0 \end{aligned}$$

Dual: (adapted from [4])

$$\begin{aligned} \max_{\alpha, \alpha^*} \quad & -\frac{1}{2} \sum_{ij} (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) \mathbf{x}_i^\top \mathbf{x}_j \\ & - \epsilon \sum_i (\alpha_i + \alpha_i^*) + \sum_i y_i (\alpha_i - \alpha_i^*) \\ \text{subject to} \quad & \sum_i (\alpha_i - \alpha_i^*) = 0 \quad \text{and} \quad \alpha_i, \alpha_i^* \in [0, C] \end{aligned}$$

Like for the one-class SVM, the weights and bias of the primal problem can be recovered from the dual solution and the KKT conditions (cf. [4]).



# Support Vector Regression

## Advantages:

- ▶ Enables the model to be robust to outliers in the data and also to small noise perturbations (→ matches well data encountered in practice).
- ▶ The SVR primal and dual both take the form of a quadratic optimization problem with linear constraints (→ one can use standard solvers and always finds the true optimum of the objective).

## Disadvantages:

- ▶ Unlike least squares regression, SVR has no closed form solution.
- ▶ Unlike least squares regression, there is no simple relation between the SVR objective and statistical measures such as correlation or variance.

# Summary

# Summary

- ▶ Regression is a very common data analysis, that gives us insights into the relation between a set of input variables and an output variable of interest, e.g. what correlation can we achieve between these variables, how the output variable responds to the different input variables.
- ▶ The most common formulation is the least square regression. It has an analytical solution and connections to the more general canonical correlation analysis.
- ▶ Least square error is not robust to outliers. In presence of outliers, it is better to consider absolute deviations. This can be addressed within the framework of support vector regression.

# Further Topics

## Ridge Regression:

- ▶ Implement a preference for models that are flatter, even at the cost of incurring additional prediction errors on the available data.
- ▶ This can be useful for improving the reproducibility of insights extracted from the model (→ Lecture 8).

## Lasso Regression:

- ▶ Implement a preference for models that respond only to a few input variables.
- ▶ This can be useful for improving the interpretability of the model and its predictions (→ Lecture 8)

## Nonlinear Regression:

- ▶ Replace the linear mapping between input and output by a nonlinear mapping, e.g. quadratic discriminants, kernel ridge regression, kernel SVR, deep neural networks.
- ▶ This enhances the representation power of regression models and enables reducing the prediction error / strengthening correlations (→ Lectures 10-12).

# References



S. Agarwal, N. Tosi, P. Kessel, S. Padovan, D. Breuer, and G. Montavon.  
Toward constraining mars' thermal evolution using machine learning.  
*Earth and Space Science*, 8(4), Apr. 2021.



H. Drucker, C. J. C. Burges, L. Kaufman, A. J. Smola, and V. Vapnik.  
Support vector regression machines.  
In *NIPS*, pages 155–161. MIT Press, 1996.



K. Preuer, G. Klambauer, F. Rippmann, S. Hochreiter, and T. Unterthiner.  
Interpretable deep learning in drug discovery.  
In *Explainable AI*, volume 11700 of *Lecture Notes in Computer Science*, pages 331–345. Springer, 2019.



A. J. Smola and B. Schölkopf.  
A tutorial on support vector regression.  
*Stat. Comput.*, 14(3):199–222, 2004.