

Thesis

Emilio Dorigatti

March 29, 2018

Acknowledgments

todo

Summary

todo

Contents

1	Introduction	7
1.1	Problem	7
1.2	Research Question	8
1.3	Research Methodology	8
1.4	Ethics and Sustainability	9
1.5	Outline	9
2	Background	11
2.1	Fluid Dynamics	11
2.1.1	Laminar and Turbulent flow	11
2.1.2	The Boundary Layer	12
2.2	The Atmospheric Boundary Layer	13
2.2.1	Surface Fluxes	14
2.2.2	The Turbulence Kinetic Energy Budget	15
2.2.3	Monin-Obukhov Similarity Theory	17
2.3	Machine Learning	18
2.3.1	Learning Theory	19
2.3.2	Cross Validation	21
2.3.3	Parameter Estimation for Regression	22
2.3.4	Hyper-parameter Optimization	24
2.3.5	Ridge Regression	25
2.3.6	k-Nearest Neighbors	25
2.3.7	Random Forests	25
2.3.8	Gaussian Processes	26
2.3.9	Neural Networks	26
3	Method	27
3.1	Data Collection	27
3.1.1	Wind Measurement	28
3.1.2	Air Temperature Measurement	28
3.1.3	Eddy Correlation	28
3.1.4	Gap Filling	29
3.1.5	Data Filtering	29
3.2	Flux-Profile Relationships	30
3.2.1	Obukhov Length	30

3.2.2	Gradients	30
3.3	Monin-Obukhov Similarity Theory	32
3.4	Model Fitting	33
3.4.1	Features	33
3.4.2	Models	34
3.5	Performance Evaluation	35
3.6	Success Criteria	35
4	Results	37
4.1	Exploratory Data Analysis	37
4.2	Gradient Computation	37
4.3	Monin-Obukhov Similarity Theory	37
4.4	Model Comparison	37
5	Discussion	41
5.1	Limitations and Future Work	41

Chapter 1

Introduction

1.1 Problem

Climatology, or climate science, studies the long-term average weather conditions. In the last few decades, climate scientists found solid evidence of ongoing changes in Earth's climate, most notably, a general increase in average temperature; in the long run, this and other changes can potentially have a devastating impact on life on our planet. Regardless of the causal effect of human activities on it, having a solid understanding of the climate allows us to find the best way to mitigate its change, and possibly prevent it completely.

Climate science is an extremely difficult field, for a number of reasons. First, climate is evident, by definition, only over long periods of time and large distances, making the usual scientific approach of testing hypotheses by conducting experiments inapplicable; instead, climate scientists have to rely on historical data. Second, the atmosphere is a complex and chaotic system, which must be described through systems of nonlinear differential equation. They can be used to create numerical simulations, and the resulting predictions compared to historical data to assess the accuracy of the theory. Furthermore, the chaotic character of the atmosphere makes it impossible to study parts of it in isolation from others, as small scale phenomena affect large scale ones, and vice versa. In spite of this, it is useful to split the atmosphere in vertical layers and horizontal zones, in order to differentiate among conditions and phenomena typically occurring in one area or the other.

Usually, the troposphere is referred to as the lowest layer of the atmosphere, but it can actually be subdivided in more sub-layers, the lowest of which is the *atmospheric boundary layer*: it is the region of the atmosphere that is affected by the conditions of the surface. Most human activities happen in this layer, and it is responsible for a large part of the diffusion and transport of aerosol such as, among others, pollutants. Yet, the physics governing the atmospheric boundary layer is not fully understood, and the theory is lacking. One important issues in the study of the atmospheric boundary layer is the derivation of flux-profile relationships for wind and temperature: they essentially relate the transport of momentum and heat by the wind (flux) with the change of wind speed/temperature with altitude (profile). The state of the art relationships are defined by the Monin-Obukhov Similarity theory in terms of the instability parameter

ξ , computed as the height above surface scaled by turbulence due to horizontal wind and vertical air movement due to variations in heat. Difficulties in measurements of relevant quantities make this theory accurate only up to 10-20%, and applicable in a restricted set of conditions.

In stark contrast to the traditional, top-down approach of science, recent developments in information technology made bottom-up approaches possible. In this new way of thinking, existing data is used to automatically infer the "best" explanation for the measurements at hand, the underlying laws that originated that data. The field that makes this possible is called Machine Learning: it takes advantage of several methods coming from statistics, information theory, optimization theory, etc., to make computers learn from examples. Together with Natural Language Processing and Automated Planning, it is one of the three main branches of Artificial Intelligence, the sub-field of Computer Science that studies ways of making machines behave intelligently.

todo more fluff about ml ?

1.2 Research Question

Currently, limitations of the validity of the Monin-Obukhov similarity theory are not believed to be a likely explanation for the high scatter that is found in experimental studies, unless in highly stable conditions Högström [1996]. With the availability of micro-meteorological data from specialized observation sites such as Cabauw, in the Netherlands, and the recent developments in Machine Learning, this conjecture can be finally put to the test. More specifically, we pose the following

Research Question 1: Is it possible to use the data from the Cesar database to improve the Monin-Obukhov model of the flux-profile relationships, by using more predictors besides the instability parameter?

Research Question 2: (in case of an affirmative answer to the first research question): What impact do the different features have on the quality of the prediction?

Affirmative answers to these question would contribute to improve the quality of current global circulation models (large scale climate simulations). Optis et al. [2014] shows that, at least in difficult conditions, several simulation models produce flux estimates that are highly inaccurate, presenting low correlation (mostly below 0.3) and errors that are as large as the fluxes themselves. Since errors in the simulation accumulate over time, even slight improvements can greatly improve the accuracy of the results after years of simulated time. Improved flux-profile relationship can, therefore, yield more accurate estimates of the fluxes.

1.3 Research Methodology

todo

1.4 Ethics and Sustainability

reproducibility: using open data, following standards of reproducible research (open source and jupyter notebooks)

1.5 Outline

todo

Chapter 2

Background

This chapter introduces the basic concepts the reader should be qualitatively familiar with, in order to understand the content of this thesis. It is assumed readers are already knowledgeable of simple mathematical concepts, such as calculus, linear algebra, statistics, and probability theory. Readers acquainted with the material should feel free to skip this chapter.

2.1 Fluid Dynamics

Fluid dynamics is the discipline that studies the flow of fluids; it has several branches that study different fluids, such as aerodynamics (the study of air motion) and hydrodynamics (the study of water motion). These disciplines are routinely applied when designing cars, airplanes, ships, pipelines, etc.

2.1.1 Laminar and Turbulent flow

There are two distinct ways in which particles in a fluid can move: laminar flow and turbulent flow. In the former, all the particles move orderly, perhaps with a different speed, but all in the same direction, whereas in the latter the movement of particles is highly chaotic and unpredictable, and tends to give rise to eddies of varying sizes. People are most familiar with the distinction between the two through the smoke rising from a cigarette, which starts smooth and becomes turbulent shortly thereafter, as in figure 2.1. The kind of flow in under specific conditions can be predicted using the Reynolds number Re , which is the ratio between inertia forces, favoring turbulent flow, and viscosity forces, stabilizing the fluid towards laminar motion:

$$Re = \frac{\rho u L}{\mu} = \frac{u L}{\nu}$$

With ρ the density of the fluid, u its velocity, L a characteristic linear dimension of the system under consideration, μ and ν the kinematic and dynamic viscosity of the fluid. The viscosity describes, intuitively, how much the molecules of the fluid tend to stick together and resist motion by generating drag. For example, water has low viscosity, and honey has high viscosity.

Figure 2.1: Smoke from a cigarette, and the transition from laminar to turbulent flow.



Since turbulence is random, it is usually studied in terms of the statistical properties of physical quantities through the Reynolds decomposition; given a quantity $a(s, t)$ which varies in space and time, we can compute its average

$$\bar{a}(s) = \frac{1}{T} \int_{T_0}^{T_0+T} a(s, t) dt$$

and the deviation from the average

$$a'(s, t) = a(s, t) - \bar{a}(s)$$

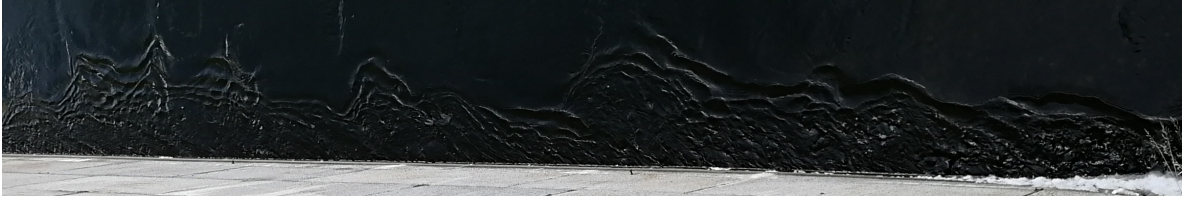
By definition, $\overline{a'} = 0$, which means that all the effects of turbulence are contained in a' . Common statistical properties such as variance and covariance are expressed respectively as $\overline{a'a'}$ and $\overline{a'b'}$.

2.1.2 The Boundary Layer

In the context of fluid dynamics, the boundary layer studies the behavior of a fluid when it is flowing close to a solid surface. Imagine a laminar flow close to a solid surface; because of viscosity, the molecules flowing near the surface move slower, and, in the limit, the velocity of the molecules in direct contact with the surface is 0 (this is called the *non-slip condition*). Thus, the velocity of the fluid increases smoothly, continuously and monotonously with the distance from the solid, until it reaches the *free-flow* velocity, after which it stays constant. The region close to the surface, where the fluid moves slower, is called the *boundary layer*, and is the region where the viscosity of the fluid influences its motion. Its height δ can be defined when the local velocity surpasses a certain threshold, such as 99% of the free-flow velocity.

The variation of velocity with distance from the surface, $\partial \bar{u} / \partial z$, is called *shear*, and, together with viscosity, determines the materialization of turbulence in the flow. Every layer of fluid is squeezed between a faster moving layer above and a slower moving below; in high shear conditions, this causes high stress on the particles, and prevents them from moving orderly, thus leading to turbulent motion. Figure 3.1 shows turbulence forming

Figure 2.2: Turbulent boundary layer at the edge of a canal; the end of the BL is clearly visible, where the flow transitions from turbulent to laminar. Water flows from right to left.



close to the wall in a canal, where the water flows from right to left. Viscosity and the no-slip condition prevent this phenomenon to arise in a region very close to the solid surface, called the *laminar (or viscous) sub-layer*, where we still find laminar motion.

The strength of the turbulence is proportional to $u_{rms} = (\overline{u'^2})^{1/2}$, which is, in turn, proportional to the shear. Again, because of the no-slip condition, u_{rms} is zero at $z = 0$, increases in the laminar sub-layer, and decreases to 0 at the end of the boundary layer, assuming laminar flow outside of it. Higher free-stream velocity generates higher shear, more turbulence, and a thinner laminar sub-layer. The strength of turbulence can be written in units of velocity, resulting in the *friction velocity*, computed as $u_* = (\tau/\rho)^{1/2} = (\nu \cdot \partial \bar{u} / \partial z)^{1/2}$, where τ is the shear stress, ρ is the density of the fluid, $\nu = \mu/\rho$ is the kinematic viscosity, and μ the dynamic viscosity. Therefore, the friction velocity increases with shear and viscosity, and decreases with density; it is proportional to the free-stream velocity and the turbulence strength, and inversely proportional to the height of the laminar sub-layer.

The mean velocity \bar{u} increases linearly within the laminar sub-layer, then logarithmically until the end of the boundary layer, thus the shear decreases further away from the surface. In the logarithmic sub-layer, the velocity is computed as $\bar{u}(z) = u_*(\log z - \log z_0)/\kappa$, where z_0 is the characteristic roughness of the surface, and κ is the von Karman's constant, whose value is around 0.4 [citation needed]. The characteristic roughness depends on the texture of the surface, and its relationship with the height δ_s of the laminar sub-layer; if the roughness scale is smaller than δ_s , the logarithmic velocity profile is not affected by the texture, because the laminar sub-layer completely covers the variations on the surface, and we have the so-called smooth turbulent flow. If, on the contrary, the bumps in the surface are larger than δ_s , the laminar sub-layer follow the profile of the surface, and the logarithmic velocity profile is altered depending on the texture, a regime called rough turbulent flow.

2.2 The Atmospheric Boundary Layer

The atmosphere is composed by air, which is behaves like fluid. Therefore, close to the Earth's surface, in the region called *atmospheric boundary layer*, we find the same effects described in the previous section. Additionally, there are other phenomena that complicate things further, such as the temperature of the surface, which changes widely from day to night and from Summer to Winter, the rotation of the Earth, the varying

roughness of the surface, due to cities and mountains, etc. The effect of the surface on the first few hundred meters of the atmosphere is the main focus of *boundary layer meteorology*.

The height of the atmospheric boundary layer (hereafter abbreviated as ABL) typically varies between 100 and 1000 meters, highly depending on the conditions, and it is always turbulent. There are three main instabilities driving turbulence in the ABL:

- Shear instability: caused by shear, the mechanism described in the previous section. This happens at high Reynolds number, and, by using typical values for the ABL, we find Re well above 10^6 .
- Kelvin-Helmholtz instability: occurs when there is a difference of density and velocity in different layers of flow. This is the mechanism that generates, for example, waves in ponds, lakes, and oceans.
- Rayleigh-Bernard instability: is caused by the decrease of potential density¹ with height, or, in other words, when warm fluid is below cold fluid; the warm fluid will rise, and the cold fluid will drop, a phenomenon called *convection*. During hot Summer days, the surface is much warmer than the air, thus the air close to the surface will heat and tend to rise.

Turbulence has the very important role of transport and mix of air properties, such as heat, moisture, particles, aerosols, etc. This is especially true in *unstable* conditions, when the air moving upwards (e.g. because it is warmer) is less dense than the air moving downwards; when the contrary happens, the ABL is called *stable*.

The ABL can be divided in two main sub-layers: the inner surface layer and the outer Ekman layer. This distinction is mainly done based on the scale of the dominating turbulent eddies: they are much smaller than the height of the ABL in the surface layer, and of comparable size in the outer layer.

It is very important to have a macroscopic understanding of the turbulent processes in the ABL, because they happen at length and time scales too small to be simulated in global climate models. The process of expressing the result of turbulence as a function of large scale parameters is called parametrization; having realistic models is essential in order to conduct precise simulations of the global climate in the scale of tens or hundreds of years, because errors tend to accumulate and amplify as the simulation goes on. Other fields that benefit from the study of the ABL are urban meteorology (interested in the dispersion of pollutants), agricultural meteorology (interested in the formation of frost and dew, the temperature of the soil, etc.), aviation (predict fog and strong winds), and so on.

2.2.1 Surface Fluxes

A flux measures the amount of a physical quantity that flows through a surface. In the context of boundary layer meteorology, we are interested in the flows through the

¹the potential density is the density that a parcel of air would attain if brought at a standard reference pressure adiabatically, i.e. disallowing exchanges of heat with its surroundings. Potential density is useful to compare densities irrespectively of pressure, i.e. altitude

surface of earth, because, through them, the surface and the atmosphere exchange energy; these fluxes are thus measured in W m^{-2} . The main source of energy for the surface is long-wave radiation coming from the sun, and short-wave radiation coming from the atmosphere and the clouds. A small amount of long-wave radiation is emitted from the surface, therefore let the net radiative flux be R , positive when the surface gains energy.

The main fluxes by which the surface loses energy to the atmosphere are called the turbulent flux of *sensible heat* H , also called kinematic heat flux, and the turbulent flux of *latent heat* λE , also called kinematic flux of water vapor/moisture. The difference between the two is that the former causes an actual change of temperature, whereas the latter does not affect temperature². The main causes of sensible heat fluxes are conduction and convection, whereas the main cause of latent heat fluxes is water movement: condensation, evaporation, melting, etc.

The final flux of interest is the soil heat flux G , which is simply the heat "absorbed" by the surface and not given to the atmosphere. These four fluxes are linked by the surface energy balance equation:

$$R = H + \lambda E + G$$

which states that the total incoming energy R must be equal to the energy given back to the atmosphere $H + \lambda E$ (not counting long-wave radiation, which is accounted to in R) plus the energy absorbed by the surface G .

The turbulent fluxes H and λE are constant in the surface layer [todo add citation]. Experimentally, the energy balance is not always achieved [todo cite from cabauw] due to difficulty in measuring fluxes due to eddy correlation being inaccurate (verify this)

2.2.2 The Turbulence Kinetic Energy Budget

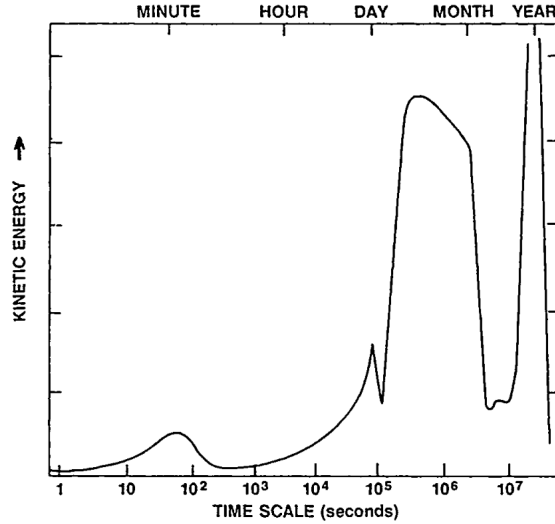
Kinetic energy is energy stored in form of movement: faster or heavier objects have more kinetic energy than slower or lighter ones. The Reynolds decomposition allows us to decompose the kinetic energy of turbulent flows in two terms: one caused by the mean flow, and one caused by turbulence. This decomposition can be justified by examining the temporal spectrum of kinetic energy, shown in figure 2.3. Four peaks are visible, corresponding to different sources of kinetic energy: turbulence, day-night cycle, westerlies³, and seasons. Importantly, there are few sources of kinetic energy in the 30 minutes to one hour time scale; this so-called spectral gap allows us to separate between turbulence and other sources of fluctuations in the atmosphere.

From now on, we will use a coordinate system with the x axis aligned to the average horizontal wind direction, the y axis perpendicular to it, and the z axis pointing away from the surface. Then, we will use the letters u , v and w to denote the components of the wind along the axes x , y and z respectively; clearly, $\bar{v} = 0$. Eddy fluxes can then

²imagine a pot of boiling water; selecting a higher temperature on the stove will not increase the temperature of water above 100°C , but will make it boil faster. The additional heat introduced in the system is dissipated through increased evaporation

³winds blowing from the east towards the west in the mid-latitudes

Figure 2.3: Change of atmospheric kinetic energy at different time-scales. The peaks in the scale of days and months and years are due to the day-night and Summer-Winter cycles, the peaks in the monthly scale are due to baroclinic instability in the mid-latitude westerlies, and the peaks at one minute are due to convection and atmospheric turbulence [Scorer, 1992, Vinnichenko, 1970]



be described in terms of covariances: let θ denote the potential temperature⁴, then $\overline{w'\theta'}$ is the turbulent heat flux, i.e. the sensible heat flux in the vertical due to wind. Usually the ABL studied assuming homogeneous horizontal conditions, because they vary on a length scale larger than the height of the ABL. Because of this, the horizontal eddy correlations $\partial \overline{u'a'}/\partial x$ and $\partial \overline{v'a'}/\partial y$ are usually of negligible intensity, and are thus ignored. Note that this is not necessarily true when clouds are involved.

It is important to notice that turbulence is dissipative in nature. Consider a hot Summer day, where air is warmer close to the surface, and a circular eddy moving some air up and some down, so that the average motion is zero. The parcel of air moving up ($w' > 0$) ends up being warmer than its surroundings ($\theta' > 0$), while the one moving down ($w' < 0$) will be colder ($\theta' < 0$); the result is a new transport of heat through the eddy: $\overline{w'\theta'} > 0$. On the contrary, imagine a cold night, where the air close to the surface is colder; the same eddy would transport a colder parcel of air upwards, and a warmer one downwards. In both cases, the end result would be a net transport of heat without transport of mass. Because of the ??? law, the eddy must lose energy, and thus dissipate over time.

Since turbulence changes over time, we are more interested in the change of kinetic energy, the *turbulent kinetic energy budget*. A full derivation is out of the scope of this work, but its final form [Högström, 1996] can be derived from prime physical principles, resulting in

⁴the potential temperature is final temperature after bringing a parcel of air to a standard pressure adiabatically, i.e. not allowing exchange of temperature with the surroundings. It is a useful mean to compare temperatures irrespectively of pressure, i.e. altitude

$$\frac{\partial \overline{e'^2}}{\partial t} = \underbrace{\overline{u'w'} \frac{\partial \bar{u}}{\partial z}}_P - \underbrace{\frac{g}{T} \overline{w'\theta'}}_B + \underbrace{\frac{\partial}{\partial z} \frac{\overline{w'e'^2}}{2}}_{T_t} + \underbrace{\frac{1}{\rho} \frac{\partial \overline{p'w'}}{\partial z}}_{T_p} + \epsilon \quad (2.1)$$

Where $e'^2 = u'^2 + v'^2 + w'^2$. The P term is the production due to shear, B is the production due to buoyancy, T_t is the turbulent transport of TKE by large-scale eddies, T_p is the transport due to pressure, and ϵ is molecular dissipation due to viscosity. P and B are the most prominent terms, and the transport terms are close to zero in neutral conditions [Högström, 1996].

The P term is always positive, as turbulence can only introduce more kinetic energy. On the other hand, the contribution from buoyancy can be either positive or negative, depending on the difference of temperature between a parcel of air moved by the turbulence and the surrounding air. When $\overline{w'\theta'}$ is negative, the turbulence is moving cold air upwards and warm air downwards; these parcels of air will try to undo the effect of turbulence, thereby increasing the overall kinetic energy. A similar reasoning goes for when the heat flux is positive.

2.2.3 Monin-Obukhov Similarity Theory

One of the factors to distinguish laminar from turbulent flow is the length scale L of the system. This length scale for the ABL was derived by A. M. Obukhov in 1946, and forms the basis of the similarity theory. According to this theory, the normalized wind and temperature profiles can be expressed as a unique function of $\xi = z/L$:

$$L = -\frac{u_*^3}{\kappa \frac{g}{\theta_v} \frac{Q}{\rho c_p}} = -\frac{u_*^3 T_v}{\kappa g w' \theta_v} \quad (2.2)$$

$$\frac{\partial \bar{u}}{\partial z} \frac{kz}{u_*} = \phi_m(\xi) \quad (2.3)$$

$$\frac{\partial \bar{\theta}_v}{\partial z} \frac{kz}{T_*} = \phi_h(\xi) \quad (2.4)$$

With

- $g = 9.81 \text{ m s}^{-2}$ the acceleration due to Earth's gravity
- $\kappa = 0.4$ the von Karman constant
- θ_v virtual temperature⁵, obtained as

$$\theta_v = \theta \frac{1 + r_v/\epsilon}{1 + r_v} = \theta(1 + 0.61 \cdot q) \quad (2.5)$$

⁵potential temperature of dry air if it had the same density as moist air. It allows to use formulas for dry air when the air is not dry.

Where θ is the air temperature, r_v is the mixing ratio, $q = r_v/(1 + r_v)$ the specific humidity, and ϵ is the ratio of the gas constants of dry air and water vapor, roughly 0.622.

- ρ the air density, computed from the pressure P and the specific gas constant for dry air $R = 287.058 \text{ J kg}^{-1} \text{ K}$ as

$$\rho = \frac{P_0}{RT_v}$$

- c_p specific heat of dry air, $1005 \text{ J kg}^{-1} \text{ K}^{-1}$ at 300 K
- Q the buoyancy flux, approximated by $H + 0.07\lambda E$ and measured in W m^{-2}
- $\overline{w'\theta_v} = Q/\rho c_p$ the flux of virtual potential temperature, measured in K m s^{-1}
- $T_* = -\overline{w'\theta}/u_*$

The stability parameter ξ is positive for stable conditions, where wind shear dominates the production of TKE, and negative for unstable conditions, where buoyancy is the main contributor to turbulence. It approaches 0 in the limit of neutral stratification (i.e. $\partial\bar{\theta}/\partial z = 0$), because the temperature flux goes to 0 causing L go to infinity.

The universal functions ϕ_m and ϕ_h must be determined experimentally. This is no easy task, and considerable effort has been devoted to it; one of the greatest difficulties lies in obtaining accurate and unbiased measurements, especially the fluxes. Högström [1988] is a meta-study that aggregates and improves many previous results, and suggests the following expressions:

$$\phi_m(\xi) = \begin{cases} (1 - 19.3\xi)^{-1/4} & -2 < \xi < 0 \\ 1 + 6\xi & 0 < \xi < 1 \end{cases} \quad (2.6)$$

$$\phi_h(\xi) = \begin{cases} 0.95(1 - 11.6\xi)^{-1/2} & -2 < \xi < 0 \\ 0.95 + 8\xi & 0 < \xi < 1 \end{cases} \quad (2.7)$$

The Monin-Obukhov similarity theory is only applicable in the surface layer, at heights much larger than the aerodynamic roughness length, and with $|\xi| < 2$; even under ideal conditions, the predictions of this theory are accurate up to 10-20% [foken 2006].

2.3 Machine Learning

The goal of Machine learning is to develop algorithms that allow computers to learn from examples. Learning is intended as the ability of inferring general rules from the available examples, so that new, previously unseen examples can be correctly characterized. The set of samples from which the computer is supposed to learn is called the *training set*, and each sample is a sequence of numbers describing its attributes, or *features*. There are three approaches in machine learning:

- *Supervised* learning: in this setting, the examples are composed of an input and a desired output, and the goal is to build a model that can correctly predict the output given the input. There are different algorithms depending on the type of output: *regression* algorithms predict continuous output, while *classification* algorithms predict discrete output.
- *Unsupervised* learning: in this setting, no output is available. The task of the algorithm is to figure out hidden relationships between the samples in the training set, for example whether they form clusters, or there are anomalous samples, or the correlation between features of the examples.
- *Reinforcement* learning: in this setting, the computer is free to act in an environment and to observe how the environment responds to its actions. Additionally, it receives a *reward* for every action it takes, and the goal of the computer is to learn a sequence of actions that maximizes the received reward. Reinforcement learning is often applied in robotics [citation] and game playing [alphazero citation].

A supervised machine learning model uses a set of parameters to compute the output value starting from the input features. The actual parameters values are learned from the training set in a process called *training*. This process is controlled by another set of parameters called hyper-parameters, whose value can be found from the training data as well. Whereas the parameters control the relationship between input and output, hyper-parameters control the "character" of the learning algorithm, such as how eager or conservative it is in learning minute details in the features. Learning too many details can be detrimental, because some differences can be due to noise, rather than actual differences.

The next sections describe the theory of learning, a general technique to estimate the parameters of a regression model, and introduce several machine learning algorithms for regression.

Notation: Scalars are denoted in *italic*, vectors (always column) and matrices in **bold**. The training set contains N training samples, indexed by n , and each sample is a pair of feature vector $\mathbf{x}_n \in \mathbb{R}^D$ and a target value $t_n \in \mathbb{R}$. The feature vectors are grouped in the $N \times D$ matrix $\mathbf{X} = [\mathbf{x}_1 \cdots \mathbf{x}_N]^\top$ and the target values in the $N \times 1$ vector $\mathbf{t} = [t_1, \dots, t_N]^\top$. Models are parametrized by a parameter vector $\boldsymbol{\theta}$ and their output for the feature vector \mathbf{x}_n is $f_n = f(\mathbf{x}_n; \boldsymbol{\theta})$. The vector containing both parameters and hyper-parameters is denoted with $\boldsymbol{\Theta}$.

2.3.1 Learning Theory

The goal of supervised learning is to use the training examples $D = (\mathbf{X}, \mathbf{t})$, independent and identically distributed according to an unknown distribution p_{XT} , to find a good prediction rule $\hat{f} : \mathcal{X} \rightarrow \mathcal{Y}$ among a family of available functions \mathcal{F} . In most practical cases, $\mathcal{X} = \mathbb{R}^D$ and \mathcal{Y} is either \mathbb{R} for regression or a subset of \mathbb{N} for classification. The goodness of a prediction rule f is measured through a *loss function* $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ that

tells, given a target value t and a guess $y = f(x)$, how much the guess is off. The *risk* of an estimator f is simply its expected loss:

$$R(f) = \mathbb{E}_{(x,t) \sim p_{XT}} [\ell(f(x), t) | f] \quad (2.8)$$

The ideal situation is to find the estimator f^* that has the lowest risk; unfortunately this is not possible, because the distribution p_{XT} is unknown. Since the training data is a sample from this distribution, we can compute the *empirical risk* of f on this set of examples instead:

$$\hat{R}(f) = \mathbb{E}_{(x,t) \sim D} [\ell(f(x), t) | f] = \frac{1}{N} \sum_{n=1}^N \ell(f(x_n), t_n) \quad (2.9)$$

We can now use the empirical risk as a surrogate for the true risk, selecting the estimator

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \hat{R}(f) \quad (2.10)$$

as our best guess. This procedure is called *empirical risk minimization*. Note that \hat{f} is a random function, because it depends on D , which is a random variable. An interesting question to ask is how good \hat{f} actually is, or, in other words, what is its expected true risk $\mathbb{E}[R(\hat{f})]$ and how it compares to the lowest attainable risk $R(f^*)$. This latter quantity can be decomposed as

$$\mathbb{E}[R(\hat{f})] - R(f^*) = \left(\mathbb{E}[R(\hat{f})] - \inf_{f \in \mathcal{F}} R(f) \right) + \left(\inf_{f \in \mathcal{F}} R(f) - R(f^*) \right) \quad (2.11)$$

Where the first term is the *estimation* error incurred by not selecting the best possible estimator in \mathcal{F} , and the second term is the *approximation* error caused by searching a good estimator in \mathcal{F} . They usually have opposite behavior:

- the estimation error is high when \mathcal{F} is too complex for the data at hand, where complexity refers to the range of phenomena that can be accurately modeled by these functions. In other words, \mathcal{F} contains many valid explanations that are equally good on the training set (they have low empirical risk \hat{R}), but are not accurate models of the underlying phenomenon p_{XT} , i.e. they do not generalize well (they have high risk R);
- the approximation error is high when \mathcal{F} cannot adequately model the phenomenon that we are trying to describe, i.e. when it does not contain any good explanation for it.

This decomposition is often referred to as the *bias-variance decomposition*, where bias and variance refer to, respectively, approximation and estimation error.

The fundamental problem is to find a class of functions that is powerful enough to model p_{XT} , but not too powerful so as to contain too many good explanations, because we would not be able to choose. Equivalently, we want to find a good model, one

that can explain the data we have, and generalize to new examples. This problem is known as *model selection*, and is important to distinguish between model selection for *identification*, and model selection for *estimation*. In the former case, the goal is to obtain a model and its parameters, to be used on new prediction problems, whereas the goal of the latter is to obtain a realistic estimate of the performance that can be obtained on new data. This problem can be approached in two ways: by directly estimating approximation and estimation errors using hold-out sets, or by penalizing complex models in order to favor simpler explanations, even though they might have slightly higher empirical risk. The next sections describe these two approaches.

2.3.2 Cross Validation

The idea behind hold-out methods is to partition the available data in two smaller sets D_T and D_V , usually of size $2/3$ and $1/3$ of the total, and use the *training* set D_T to choose \hat{f} and D_V to estimate its risk. Since D is assumed to be a representative sample from p_{XT} , if the two partitions contain independent and identically distributed samples, the empirical risk on the *validation set* D_V can give us a glimpse on the generalization power of \hat{f} . This is because the validation set contains new samples that were not used to choose \hat{f} , thus the empirical risk on this set is an unbiased estimate of the true risk of the discovered model. This allows us not only to be confident about the performance of the estimator on unseen data, but also to compare different estimators. The problem of a single hold-out set is the variance of its estimate of the risk, which depends on the size of the validation set. This means that we are faced with a trade-off: use a lot of data to select a good estimator, but have high uncertainty in its estimated performance, or use less data and select a less powerful estimator, but have a more accurate picture of its performance.

The solution to this problem is to repeatedly perform this partitioning procedure so as to obtain many estimates of the risk, each on a different subset of validation samples, and average these results together. This can be done in a number of different ways:

- in random subsampling, the procedure above is simply repeated many times, by using two thirds random examples for training, and the remaining one third for validation;
- In bootstrapping [Efron and Tibshirani, 1993], the training set is created by taking $N = |D|$ examples *with replacement* from D , and using the remaining examples for validation. This means that the validation set contains on average approximately 36.8% of the samples in D , and the training set the remaining 63.2%, with many duplicates;
- In k -fold cross validation [Geisser, 1975], D is partitioned in k subsets, and each of them is used in turn as validation set, while the others are used for training. This produces k estimates of the true risk, coming from the k subsets.

Regardless of the method used, the final estimate of the performance is the average of the estimates obtained from the individual trials. Every method has different properties regarding both the bias and the variance of the estimates, and there is considerable

controversy on which method should be used in which situation. For example, Kohavi [1995] recommends using 10-fold cross validation for comparing models, because, although its estimate of the performance is biased, it has lower variance compared to bootstrapping; however, Bengio and Grandvalet [2004] show that it is not possible to obtain an universal unbiased estimate of the variance of k-fold cross validation. Zhang and Yang [2015] further discusses this issue, and debunks some myths and commonly held misconceptions about cross validation, including the belief, consequent Kohavi [1995], that 10-fold cross validation is always the best choice. Generally, there is a tradeoff in choosing the value of k, as high values yield estimates with lower bias, but higher variance [Arlot et al., 2010], and are more computationally intensive.

nested cv Stone [1974]

2.3.3 Parameter Estimation for Regression

In this section, we describe a general framework, rooted in Bayesian statistics, for estimating the parameters of a regression model, while controlling overfitting. An advantage of Bayesian methods is that they offer a sound theoretical foundation for model selection, without requiring repeated experiments to choose among candidate models, although this mathematical rigor is not free from practical difficulties [Wasserman, 2000, Chipman et al., 2001].

A common assumption in the regression setting is that the observations are corrupted by additive Gaussian white noise, i.e. $t_n = y_n + \epsilon_n$, where y_n is the "true" value, and $\epsilon_n \sim \mathcal{N}(0, \beta^{-1})$ is the noise. Let $f_n = f(\mathbf{x}_n; \boldsymbol{\theta})$ be the model's prediction for the sample \mathbf{x}_n , then we can write the probability of observing t_n , assuming that $f_n = y_n$, as:

$$p(t_n | \mathbf{x}_n, \boldsymbol{\Theta}) = \mathcal{N}(t_n | f(\mathbf{x}_n; \boldsymbol{\theta}), \beta^{-1}) \quad (2.12)$$

This probability is called *likelihood* of the observation t_n , under the model $f(\cdot; \cdot)$ with parameters $\boldsymbol{\theta}$. Since the training data is assumed to be independent and identically distributed, the likelihood of the whole training set is

$$p(\mathbf{t} | \mathbf{X}, \boldsymbol{\Theta}) = \prod_{n=1}^N \mathcal{N}(t_n | f(\mathbf{x}_n; \boldsymbol{\theta}), \beta^{-1}) \quad (2.13)$$

And the predicted value t for a new sample \mathbf{x} is distributed as

$$p(t | \mathbf{x}, \mathbf{X}, \mathbf{t}) = \int p(t | \mathbf{x}, \boldsymbol{\Theta}) \cdot p(\boldsymbol{\Theta} | \mathbf{t}, \mathbf{X}) \, d\boldsymbol{\Theta} \quad (2.14)$$

In practice, it is not feasible to compute this integral, and its value is dominated by the values of $\boldsymbol{\theta}$ close to the one that maximizes equation 2.13 anyways; this is the gist of *maximum likelihood estimation* (MLE). Note that, since the goal is to predict y_n , there is a high probability that the "true" parameters would *not* be the ones with maximum likelihood. When a model learns the noise in the training data, it cannot generalize well to new, unseen data, because the noise is random. This situation is known as *overfitting*,

and tends to happen when the model is too complex relative to the amount of data available for training, and is related to high estimation error mentioned previously.

The risk of overfitting can be reduced with a number of *regularization* strategies. A widely used strategy consists in including a prior distribution on the parameters of the model, and maximizing their posterior distribution, computed using Bayes theorem:

$$\begin{aligned} p(\Theta|\mathbf{t}, \mathbf{X}) &= \frac{p(\mathbf{X}, \mathbf{t}|\Theta) \cdot p(\Theta)}{p(\mathbf{X}, \mathbf{t})} \\ &\propto p(\mathbf{t}|\mathbf{X}, \Theta) \cdot p(\mathbf{X}|\Theta) \cdot p(\Theta) \\ &\propto p(\mathbf{t}|\mathbf{X}, \Theta) \cdot p(\Theta) \end{aligned} \quad (2.15)$$

where we removed $p(\mathbf{X}, \mathbf{t})$ and $p(\mathbf{X}|\Theta) = p(\mathbf{X})$ because they are constant for a given dataset, and we are not interested in the exact probability, but where it reaches its maximum value. This parameter estimation procedure is called *maximum a posteriori* (MAP). Two commonly used prior distributions are the multivariate Laplace and the multivariate Normal, leading respectively to L1 and L2 regularization, when centered and symmetrical/spherical.

The maximization of the posterior can be done conveniently by maximizing its logarithm; this gives expressions that are easier to handle analytically, and give less numerical problems when computed. The MAP problem can be formulated as follows:

$$\Theta^* = \underset{\Theta}{\operatorname{argmax}} \log p(\mathbf{t}|\mathbf{X}, \Theta) + \log p(\Theta) \quad (2.16)$$

If we assume a Gaussian likelihood like the one in equation 2.13, and a spherical Gaussian prior distribution $p(\Theta) = \mathcal{N}(\mathbf{0}, \lambda \mathbf{I})$, the MAP estimation of equation 2.15 becomes, after removing unnecessary constant terms,

$$\theta^* = \underset{\theta}{\operatorname{argmin}} L(\theta) = \underset{\theta}{\operatorname{argmin}} \sum_{n=1}^N (f(\mathbf{x}_n, \theta) - t_n)^2 + \lambda \theta^\top \theta \quad (2.17)$$

Where L is the *loss function*, in this case the sum-of-squares error. It is customary to use the mean squared error instead of the sum of squares because it results in smaller numbers and is readily interpreted; being only a constant factor away from 2.17, it does not transcend the essence of MAP estimation. The parameter β can be estimated as well in a similar way, if necessary, and a fully Bayesian treatment allows to estimate λ , too. Depending on the model $f(\mathbf{x}_n, \theta)$, 2.17 can be solved analytically to yield a closed-form solution for θ .

When this is not possible, iterative optimization methods are employed. A widely used approach is called *gradient descent*: the gradient of a function computed at a given location "points" to the steepest direction where the function's value increases. By repeatedly following the gradient, it is possible to reach a local maximum, and, in the opposite direction, a local minimum:

$$\theta_{n+1} := \theta_n - \eta \cdot \nabla f(\mathbf{x}_n) \quad (2.18)$$

the series $\theta_1, \dots, \theta_n$ is guaranteed to converge to the local optimum at a rate of $O(n^{-1})$ if certain conditions are met [Gitman et al., 2018]. Nocedal and Wright [1999] describes more advanced optimization techniques that use the gradient and, possibly, the Hessian, such as Newton's. The gradient of $L(\theta)$ is

$$\nabla_{\theta} \log p(\theta | \mathbf{t}, \mathbf{X}) = 2 \sum_{n=1}^N (f(\mathbf{x}_n; \theta) - t_n) \cdot \nabla_{\theta} f(\mathbf{x}_n; \theta) + 2\lambda \theta \quad (2.19)$$

And its Hessian is

$$\nabla_{\theta}^2 \log p(\theta | \mathbf{t}, \mathbf{X}) = 2 \sum_{n=1}^N \left[(\nabla_{\theta} f(\mathbf{x}_n; \theta)) (\nabla_{\theta} f(\mathbf{x}_n; \theta))^{\top} + (f(\mathbf{x}_n; \theta) - t_n) \nabla_{\theta}^2 f(\mathbf{x}_n; \theta) \right] + 2\lambda \mathbf{I} \quad (2.20)$$

In some cases, the gradient and the Hessian can be approximated using a random subset of the training set, and, in extreme cases, a single sample. These variants are called *minibatch* gradient descent and *stochastic* gradient descent respectively. They both compute a noisy approximation to the true gradient, which can actually improve convergence and generalization of high-dimensional, non-convex loss functions such as those found in deep learning [Neelakantan et al., 2015, Smith and Le, 2017]. Vanilla gradient descent can be greatly improved with a number of techniques, such as momentum [Rumelhart et al., 1986], adaptive learning rate [Duchi et al., 2011, Zeiler, 2012, Kingma and Ba, 2014], and so on, see Ruder [2016] for an overview.

2.3.4 Hyper-parameter Optimization

The previous section discussed some methods to find the best parameters for a model. In practice, though, finding good values for the hyper-parameters is often as important as fitting the model, since these hyper-parameters control the learning process itself.

A simple way to approach this problem is to choose some possible values for each hyper-parameter, try all the possible combinations, and pick the one that works best. Alternatively, one can sample each hyper-parameter from a probability distribution, and try many random combinations of values; Bergstra and Bengio [2012] showed that this latter method is surprisingly effective at this task, and scales much better than the former. Unfortunately, this procedure can overfit, too, so it has to be paired with some resampling technique such as cross validation. The procedure is simply to test every hyper-parameter combination on every fold, so that 50 combinations tested with 10-fold cross validation require to fit the model 500 times.

More sophisticated techniques treat hyper-parameter optimization as a regression problem, and use supervised machine learning to predict the performance of a given hyper-parameter combination, and to guide the search accordingly. Clearly, one should use models that are not sensitive to their own hyper-parameter setting, or the problem is just moved! For this reason, popular algorithms that are used for this are evolutionary algorithms [citation] and Bayesian non-parametric models [citation].

2.3.5 Ridge Regression

A linear regression model has the form $f(\mathbf{x}_n; \boldsymbol{\theta}) = \boldsymbol{\theta}^\top \mathbf{x}_n$; ridge regression is simply L2-regularized linear regression. This model is simple enough that the solution for equation 2.17 can be found analytically in closed form:

$$\boldsymbol{\theta}^* = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{t} \quad (2.21)$$

Notice that the term *linear* in linear regression refers to linearity with respect to the parameters, not the features. In fact, new features can be added through the use of a possibly non-linear feature mapping $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$, such that $\mathbf{x}'_n = \phi(\mathbf{x}_n) = [\phi_1(\mathbf{x}_n), \dots, \phi_M(\mathbf{x}_n)]^\top$. Then, the model parameters can be fitted to the augmented training set \mathbf{X}' . A typical feature added is a bias $\phi_1(\mathbf{x}) = 1$.

2.3.6 k-Nearest Neighbors

The k-nearest neighbor model Stone [1977] predicts a value for a sample \mathbf{x} by averaging the target values of the K training samples that are closest to \mathbf{x} , possibly weighted by distance or some other metric. Let $\mathbf{t}_{(1)}, \dots, \mathbf{t}_{(k)}$ be the target values of the K training samples closest to \mathbf{x} according to a distance metric $d : \mathbb{R}^D \rightarrow \mathbb{R}$, then we have

$$f(\mathbf{x}) = \sum_{k=1}^K w_{(k)} \cdot \mathbf{t}_{(k)} \quad (2.22)$$

Note that the k-nearest neighbors algorithm does not have parameters, only hyper-parameters: K itself, the distance function, and the weighting scheme. Typical distance functions are the Manhattan distance $d(\mathbf{u}, \mathbf{v}) = \sum_{d=1}^D |u_d - v_d|$ and the squared Euclidean distance $d(\mathbf{u}, \mathbf{v}) = \sum_{d=1}^D (u_d - v_d)^2$, and typical weights are uniform $w_{(k)} = 1/K$ and based on distance $w_{(k)} = d(\mathbf{x}, \mathbf{x}_{(k)}) / \sum_{k'=1}^K d(\mathbf{x}, \mathbf{x}_{(k')})$.

2.3.7 Random Forests

Random forests [Breiman, 2001] are an ensemble approach to prediction whose output values are the combination of the output of multiple decision trees. In general, ensemble methods improve the performance of a single model by combining the output of many models trained on random subsets of features and/or samples.

A decision tree is a sequence of nested *if-then-else* rules, each of them testing one of the features of the input sample: inference starts from the root, and successively traverses each node by moving to the correct child, until a leaf, containing the final prediction, is reached. For classification, the prediction is the most frequent class among the training samples that end up in that leaf, while for regression the output value can be their average, or the output of a more complicated regression model trained on those samples.

Given a training set, many decision trees can be built using different random subsets of features. This creates a number of *expert* trees, each of them specialized in a restricted domain of the data. No tree is good enough on its own as it overfits its domain, but by averaging the output of many decision trees, their errors tend to cancel out.

Breiman et al. [1984] introduced a simple and widely used algorithm for building decision trees, CART; it builds decision trees by greedily and recursively partitioning the training samples until a stopping criterion is met. Partitions are created by selecting the split that has the smallest impurity; a common measure for regression is the mean squared error, while Gini impurity or cross-entropy are popular options in classification. Let X be a set of training samples and $\Sigma = (f, \tau)$ be a candidate split on feature f with threshold τ , then we can partition the training samples in two sets:

$$S^-(\Sigma) = \{\mathbf{x} \in X | \mathbf{x}_f \leq \tau\} \quad (2.23)$$

$$S^+(\Sigma) = \{\mathbf{x} \in X | \mathbf{x}_f > \tau\} \quad (2.24)$$

And define the impurity I of Σ as the weighted average of the impurities of the two sets:

$$I(\Sigma) = \frac{|S^-(\Sigma)|}{|X|} I(S^-(\Sigma)) + \frac{|S^+(\Sigma)|}{|X|} I(S^+(\Sigma)) \quad (2.25)$$

The best split is then

$$\Sigma^* = \underset{\Sigma}{\operatorname{argmin}} I(\Sigma) \quad (2.26)$$

This best split defines a new node in the tree, and the construction can proceed recursively on its two children $S^-(\Sigma^*)$ and $S^+(\Sigma^*)$. Stopping criteria for the construction can be a minimum impurity $I(X)$, a minimum impurity decrease with respect to the parent $I(X) - I(\Sigma^*)$, or a minimum cardinality for X or the two splits S^+ and S^- .

2.3.8 Gaussian Processes

todo

2.3.9 Neural Networks

todo He et al. [2015]

Chapter 3

Method

This chapter describes how the data is collected (section 3.1) and used to re-create the flux-profile relationships (section 3.2), as well as their prediction based on the Monin-Obukhov similarity theory (section 3.3). Then, we describe the main contribution of this thesis, namely how their prediction can be improved using machine learning techniques (section 3.4). Finally, we describe how the evaluation is performed (section 3.5), and explicitly state the necessary criteria for a successful answer to the research questions (section 3.6).

3.1 Data Collection

The Cabauw Experimental Site for Atmospheric Research¹ (Cesar) is a consortium formed by eight Dutch institutes and universities, which collaborate to operate and maintain an observatory for micro-meteorological conditions near the village of Cabauw, the Netherlands. The data collected characterizes the state of the atmosphere and the soil, and their interaction via radiation and surface fluxes.

The observatory is surrounded by fields and no urban agglomerations is present within 15 kilometers; the land is flat with changes of altitude within a few meters over 20 kilometers. The main mast is 213 meters high and offers measurement levels every 20 meters; at each level there are three booms of length 10 meters that allow observations unobstructed by the main mast. There are three additional smaller masts of height 10 and 20 meters located close to the main mast, in order to obtain undisturbed measurements at the lower levels, and facilities to perform soil and surface observations.

The main focus of this project is on wind and temperature profiles, and the turbulent fluxes of sensible and latent heat. Additional variables, such as temperature and humidity, are needed to compute quantities of interest, most importantly the Obukhov length, and as possible predictors for the universal functions. There is one measurement for each variable every ten minutes, and missing measurements are gap-filled with a number of techniques. The data collected is always visually validated by an operator, which marks suspect or invalid sections of data

The Cesar observatory provides full information regarding data collection², what

¹<http://www.cesar-database.nl/>

²<http://projects.knmi.nl/cabauw/insitu/observations/documentation>

Figure 3.1: Sonic anemometer to measure wind and optical open-path sensor to measure humidity at the 180m level of the Cabauw mast.



follows is a brief summary of the most relevant sections.

3.1.1 Wind Measurement

Wind speed and direction are measured at heights of 200, 140, 80, 40, 20 and 10 meters, in either two or all three booms available. The wind vane that measures direction has a resolution of 1.5° , and the cup anemometer that measures wind speed has an accuracy of the largest between 1% and 0.1 m s^{-1} . Monna [1978] studied the threshold sensitivity of both instruments, and found it lower than 0.5 m s^{-1} , even though the measurements are inaccurate up to 3 m s^{-1} .

For every ten minutes interval, the measurement comes from the instrument that is best exposed to the wind, and less affected by the obstruction caused by the mast. Corrections are then applied to the raw measurements to further attenuate the disturbance by the tower, following Wessels [1983].

3.1.2 Air Temperature Measurement

Air temperature is measured at heights of 200, 140, 80, 40, 20, 10 and 1.5 meters. todo write more

3.1.3 Eddy Correlation

The eddy correlation technique is used to compute the turbulent surface fluxes of sensible and latent heat, as well as momentum and CO_2 , starting from fluctuations in wind, temperature, and humidity.

These measurements are obtained with, respectively, a sonic anemometer, a sonic thermometer, and an optical open-path sensor. Sonic anemometers measure the wind speed by leveraging the fact that the speed of sound in free air is affected by the speed of the air itself; since the speed of sound is known, the wind speed can be easily recovered from the time a sound impulse takes to travel a short distance. By measuring the wind velocity along three orthogonal paths, the full wind vector can be recovered. Sonic thermometers work similarly, by leveraging the fact that the speed of sound is affected by the temperature of the medium it travels in. These instruments can take up to 100 measurements per second. Optical open-path sensors quantify the amount of water vapor and carbon dioxide in the air by emitting a ray of infrared light and measuring its intensity 10 to 20 centimeters further away. H_2O and CO_2 molecules in the air absorb electromagnetic radiation at known frequencies, thus the concentration of water vapor and carbon dioxide can be inferred by measuring the attenuation at these wavelengths.

The eddy correlation technique measures fluxes by computing their sample covariance with the vertical wind speed. Let w_t be the vertical wind speed at time t , then the turbulent vertical flux for the quantity a is computed as follows:

$$F_a = \frac{1}{T_2 - T_1} \sum_{t=T_1}^{T_2} (w_t - \bar{w})(a_t - \bar{a})$$

Where \bar{w} and \bar{a} are the averages of w_t and a_t for $T_1 \leq t \leq T_2$. The fluxes in the Cesar database are computed every ten minutes, with 10 measurements per second.

The eddy correlation technique is far from perfect, see Lee et al. [2004] for a detailed summary of issues.

3.1.4 Gap Filling

With gap-filling, missing measurements are replaced by synthetic values. The gap-filling method depends on the missing parameter and the duration of the period where data is not available. There are two classes of parameters: forcing parameters, which include wind, temperature, specific humidity, incoming radiation and rain, and validation parameters, which include the surface fluxes, outgoing radiation, and friction velocity.

For less than two hours of missing measurements, both forcing and evaluation parameters are gap-filled by interpolation of nearby values. For longer periods, forcing parameters are derived by transforming measures obtained from the nearby site of De Bilt, which are themselves gap-filled, if necessary. Evaluation parameters are computed with a vegetation model that uses the forcing parameters as input. The gap-filling procedure is performed by the Cesar consortium.

3.1.5 Data Filtering

Following other works in this field, such as Korrell et al. [1981] and Höglström [1988], we exclude all the records where any of the following conditions applies:

- The sensible heat flux H is smaller than 10 W m^{-2} ;

- The friction velocity u_* is smaller than 0.1 m s^{-1} ;
- The wind speed \bar{u} is lower than 1 m s^{-1} .

todo also mention dew point, soil heat (and that missing values are not imputed), net radiation, rain

3.2 Flux-Profile Relationships

Since the turbulent fluxes and the friction velocity are measured at the surface level, we can compute the flux-profile relationships only at 10, 20 and 40 meters, because these quantities can be assumed constant only within the surface layer. It is very hard to know the exact height of the surface layer, because it depends on the weather and no exact formulas are known, but it is usually assumed to be 10% as high as the boundary layer. Based on the typical height of the boundary layer, the surface layer is often higher than 40 meters and lower than 80. JW Verkaik [2006] indeed reports that a large number of observations from Cabauw at 20 m are inside the surface layer, the 100 m level is already outside of the surface layer, and Korrell et al. [1981] used the observations at 50 m in their analysis, but not those at 100 m.

3.2.1 Obukhov Length

The Obukhov length is computed as in equation 2.2, reported here:

$$L = -\frac{u_*^3}{\kappa \frac{g}{\theta_v} \frac{Q}{\rho c_p}} = -\frac{u_*^3 T_v}{\kappa g w' \theta_v}$$

The flux of virtual potential temperature can be computed following the formulas in section 2.2.3, as the data contains all the necessary quantities; u_* is given, as well as the specific humidity and the pressure. The Obukhov length is computed at each height level using the corresponding air temperature measurement, and the fluxes measured at the surface.

3.2.2 Gradients

The flux-profile relationships are listed in equations 2.3 and 2.4, and are reported here for the reader's convenience:

$$\phi_m(\xi) = \frac{\partial \bar{u}}{\partial z} \frac{kz}{u_*}$$

$$\phi_h(\xi) = \frac{\partial \bar{\theta}_v}{\partial z} \frac{kz}{T_*}$$

In order to compute them from the data, we need to compute the derivative of the wind speed, for ϕ_m , and of the virtual temperature, for ϕ_h . In general, the derivative can be obtained by fitting a model on the observations, and computing the derivative using

the model. The simplest option is to use a piecewise linear function that passes through the measurements; let the observations be sorted by height and y_i the measurement at height z_i , then for $z \in [z_i, z_{i+1}]$ we have:

$$f(z) = y_i + (z - z_i) \frac{y_{i+1} - y_i}{z_{i+1} - z_i} \quad (3.1)$$

The derivative of this function at height z_i is then the average of the slope of the segments that start and end at z_i :

$$\begin{aligned} f'(z_i) &= \lim_{h \rightarrow 0} \frac{f(z_i + h) - f(z_i - h)}{2h} \\ &= \frac{1}{2} \cdot \lim_{h \rightarrow 0} \left(\frac{f(z_i + h)}{h} - \frac{f(z_i - h)}{h} \right) \\ &= \frac{1}{2} \cdot \lim_{h \rightarrow 0} \left(\frac{y_i}{h} + \frac{y_{i+1} - y_i}{z_{i+1} - z_i} - \frac{y_i}{h} + \frac{y_i - y_{i-1}}{z_i - z_{i-1}} \right) \\ &= \frac{1}{2} \cdot \left(\frac{y_{i+1} - y_i}{z_{i+1} - z_i} + \frac{y_i - y_{i-1}}{z_i - z_{i-1}} \right) \end{aligned} \quad (3.2)$$

In order to compute this derivative at the lowest measured level $z_1 = 10$ m, we can exploit the no-slip condition and introduce an artificial observation $y_0 = 0 \text{ m s}^{-1}$ at z_0 . The value of the *roughness length* z_0 depends on the properties of the surface, and, although no unambiguous value is known for the Cabauw observatory, its value is likely between 10^{-1} m and 10^{-2} m JW Verkaik [2006], Baas et al. [2017], therefore it is reasonable to conclude that its effect is negligible on the final gradient.

A more complicated model, introduced by T. M. Nieuwstadt [1984], is

$$f(z) = a + bz + cz^2 + d \ln z \quad (3.3)$$

The model is linear in its parameters, therefore equation 2.21 can be used to compute the coefficients, using the feature mapping $\phi(z) = [1, z, z^2, \ln z]^\top$ and regularization parameter $\lambda = 0$. Once the coefficients are known, the gradient is trivial to compute:

$$f'(z) = b + 2cz + d \frac{1}{z} \quad (3.4)$$

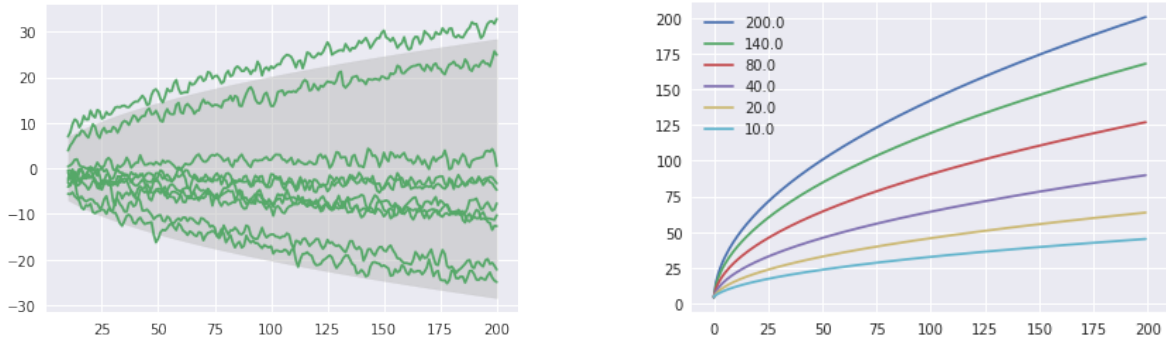
After removing the observations for which this model has a $R^2 > 0.9$, T. M. Nieuwstadt [1984] estimates the uncertainty of this gradient to be around 30%.

Finally, a third method of computing the gradient is to fit a Gaussian process to the profile using the following kernel:

$$k(z_1, z_2) = k + \exp \left(-\frac{(z_1 - z_2)^2}{2\sigma_0^2} \right) + \sqrt{\sigma_1^2 z_1 z_2 + \sigma_2^2} \mathbb{1}[z_1 = z_2] \quad (3.5)$$

The parameter σ_2 controls the noise, and depends on the instrument used to measure the modeled quantity; for example, for wind, we have $\sigma_2^2 = 0.1 \text{ m}^2 \text{ s}^{-2}$; the other parameters are found by optimizing the marginal likelihood.

Figure 3.2 shows the prior distribution of a Gaussian process with this kernel, as well as the value of the kernel for some choices of z_1 and z_2 .



(a) Prior distribution and some samples; the grey area is the 95% confidence interval.

(b) Values of the kernel: each line corresponds to a different value of z_1 , and z_2 is on the ordinate.

Figure 3.2: Behavior of the kernel in equation 3.5, where the altitude is on the x axis, and the predicted value on the y axis. Notice the approximate logarithmic profile with which the predicted value changes with altitude; this is caused by the square root term in the kernel, and emulates the effect of the $\ln z$ term in equation 3.3.

3.3 Monin-Obukhov Similarity Theory

Even though there is agreement on form of their form, there is still debate on the exact values of the coefficients, with different experiments resulting in different values [Högström, 1988]. In order to ensure a fair comparison with the results of this work, we fit the universal functions to the data from the Cesar database. Their general form is:

$$\phi(\xi) = \begin{cases} a + b\xi & \xi \geq 0 \\ a(1 - c^2\xi)^d & \xi < 0 \end{cases} \quad (3.6)$$

Where a is close to 1, b is positive, and d is negative. Since d is negative, the base of the power must be positive, hence the squared c . Following the approach outlined in section 2.3.3, the coefficients a , b , c and d can be found by minimizing the L2-regularized mean squared error using equations 2.19 and 2.20; in this case, the parameter vector is $\theta = [a, b, c, d]^T$. The gradient of 3.6 is:

$$\nabla_{\theta}\phi(\xi)|_{\xi \geq 0} = \begin{bmatrix} 1 \\ \xi \\ 0 \\ 0 \end{bmatrix} \quad \nabla_{\theta}\phi(\xi)|_{\xi < 0} = \begin{bmatrix} \tau^d \\ 0 \\ -2acd\xi\tau^{d-1} \\ a\tau^d \ln \tau \end{bmatrix} \quad (3.7)$$

where $\tau = 1 - c^2\xi$. The Hessian of 3.6 when $\xi \geq 0$ is simply 0, because it is a linear function in all parameters, while, in the negative case, we have:

$$\nabla_{\theta}^2 \phi(\xi)|_{\xi < 0} = \begin{bmatrix} 0 & 0 & -2cd\xi\tau^{d-1} & \tau^d \ln \tau \\ 0 & 0 & 0 & 0 \\ -2cd\xi\tau^{d-1} & 0 & \frac{2ad\xi\tau^d}{c^4\xi^2+\tau}(2c^2d\xi - c^2\xi - 1) & -2ac\xi\tau^{d-1}(d \ln \tau + 1) \\ \tau^d \ln \tau & 0 & -2ac\xi\tau^{d-1}(d \ln \tau + 1) & \tau^d \ln^2 \tau \end{bmatrix} \quad (3.8)$$

Analytical computation of the Hessian allows us to use the Newton conjugate gradient descent algorithm, which provides super-linear convergence rate, unlike other conjugate gradient methods whose rate of convergence is only linear [Nocedal and Wright, 1999].

This model is then fitted to the data and evaluated as the other regression models; see section 3.5 for details on the procedure.

3.4 Model Fitting

In this section, we discuss how we use the data from the Cesar database to predict ϕ_m .

3.4.1 Features

All the features come from the Cesar database, refer to section 3.1 for details on how the data was collected. The predictors are partitioned in five sets:

- F1: altitude z , wind at z , temperature at z , wind at 10 meters, temperature at 10 meters, wind at 20 meters, temperature at 20 meters, wind at 40 meters, temperature at 40 meters, soil temperature;
- F2: Soil heat flux;
- F3: Net radiation;
- F4: Rain amount, dew point;
- F5: Turbulent kinetic heat flux, turbulent latent heat flux;

These sets are used cumulatively in the order they are listed, meaning that, for example, F3 is used in conjunction with F2 and F1. Derived features that can be computed from others, such as the virtual temperature (equation 2.5), are not included. The reason is that the less inputs a model requires, the more useful and "agile" it can be. This division was decided based on domain knowledge, so that the level of a feature reflects both its expected impact on the performance and how desirable it is to include it. An example of the former reasoning is with F4, where the effect of moisture is expected to be captured by the soil heat flux, and to be generally negligible in all but the most extreme conditions. An example of the latter reasoning is with F5, because turbulent fluxes are hard to measure accurately (see section 3.1.3), and current simulation models are known to be quite inaccurate in their estimation of these fluxes [Optis et al., 2014]. Similarly, the friction velocity was not included, because the point of predicting ϕ_m is

to use it to estimate u_* from the wind gradient, which is readily measured both in real life and in simulations.

We also create a second version of each feature set, augmented with the hourly trend of each variable, except for z . The reason for including the trend is that it can give an indication of, for example, the time of the day, or other complex phenomena for which there is no measurement. The interval for the trend (one hour) is used because it is enough to capture local variations, but not too large so as to contain irrelevant information. The hourly trend is computed simply as the difference between the current value and the value measured one hour before, divided by one hour; given the goal of this work, it would not be feasible to use the future value to compute the trend.

Finally, each feature is centered and standardized so as to have zero mean and unit standard deviation:

$$z_{i,j} = \frac{x_{i,j} - \mu_j}{\sigma_j} \quad (3.9)$$

This method is not robust to outliers, since they heavily affect mean and standard deviation; this can be prevented by subtracting the median and normalizing with the interquartile range. These two methods give similar results in our datasets, therefore we follow equation 3.9, since it is the most widely used in practice. Obviously, μ_j and σ_j are computed only on the training data, and used to normalize both train and test data.

3.4.2 Models

The models that we use are ridge regression, k-nearest neighbors, gradient boosted trees, and neural networks. Due to the size of the dataset, we cannot use kernel-based algorithms such as SVM and Gaussian Processes, as they require $\Theta(N^2)$ memory to store the kernel matrix, and $O(N^3)$ time to invert it. Gradient boosted trees were used instead of random forests because they were shown to work slightly better [citation], and small-scale experiments on our dataset showed they can be fit almost an order of magnitude more quickly.

All models but neural networks are fit using nested cross validation with random hyper-parameter search. This procedure detailed in section 3.5, but we list here the distributions of the hyper-parameters that we used in the random search:

- Ridge regression: the only hyper-parameter to tune is the regularization coefficient, for which we use a \log_{10} -uniform distribution from 10^{-6} to 10_0^3 .
- k-nearest neighbors: the hyper-parameters to tune are the number of neighbors, chosen uniformly from 1 to 15, the distance function, either L1 or L2 norm, and the weights of the neighbors, either uniform or directly proportional to the distance to the query point.
- Gradient boosted trees: asd

³a random variable X has a \log_β -uniform distribution from a to b if $\log_\beta X$ is uniformly distributed between a and b . Equivalently, if Y is uniform between a and b , then β^Y is \log_β -uniform.

3.5 Performance Evaluation

The goal of performance evaluation is to obtain a realistic and unbiased estimate of the performance of a model on unseen data; we are not interested in the values of the hyper-parameters that yield this performance. For this reason, we perform nested cross validation, with the inner loop used to optimize the hyper-parameters through random search [Bergstra and Bengio, 2012], and the outer loop used to evaluate the model. We use 10 folds in both the inner and outer CV, and test 60 different hyper-parameter combinations. This number was chosen based following the advice by Bergstra and Bengio [2012], based on a simple probabilistic argument: if every combination has a 5% probability of giving a model whose performance is within 5% of the best attainable performance, then the probability that at least one in 60 trials finds a combination within 5% of the optimal is $1 - 0.95^{60} \approx 0.95$. Algorithm 1 shows in detail how nested cross validation with random search works.

A fundamental assumption underlying hold-out evaluation methods is that the samples in the training set are independent and identically distributed, so that the distribution in the two partitions are equal. This is not our case, since there is a inherent time dependency in the data, meaning that samples obtained close in time are very similar. This can be confirmed by training and evaluating a k-nearest neighbors classifier with $k = 1$ on random splits: the resulting mean squared error is in the order of 10^{-3} , clearly unrealistic. To circumvent this problem, the CV folds are created on *months*: all the samples in a given month and year are either in the training set or in the validation set.

The main evaluation metric is the mean squared error, but we compute other metrics in the outer CV to get a more complete idea of the performance of the estimators:

- Mean Squared Error
- Mean Absolute Error
- Median Absolute Error
- Mean Absolute Percent Error
- Median Absolute Percent Error
- R^2 Score:

$$1 - \frac{\sum (f_n - t_n)^2}{\sum (t_n - \bar{t})^2}$$

Where f_n is the predicted value for the test sample \mathbf{x}_n with true value t_n , the squared error is $(f_n - t_n)^2$, the absolute error is $|f_n - t_n|$, the absolute percent error is $|1 - f_n/t_n|$

3.6 Success Criteria

A successful answer to the first research question entails finding a model whose mean squared error on F5 with trend is smaller than the mean squared error of the Monin-Obukhov similarity theory estimator introduced in section 3.3. The second research question is answered by comparing the different evaluation metrics on the ten feature sets (F1 to F5, with and without trend).

Algorithm 1 Nested cross validation with random hyper-parameter search.

```

for  $k_o \leftarrow 1 \dots K_O$  do                                 $\triangleright$  Outer  $K_O$ -fold cross validation
  Generate  $k_o$ -th outer fold  $(D_T^o, D_V^o)$  from the dataset  $D$ 
   $best\_mse \leftarrow \inf$ 
   $best\_params \leftarrow \perp$ 
  for  $r \leftarrow 1 \dots R$  do                                 $\triangleright$  Find best hyper-parameters on  $D_T^o$ 
     $params \leftarrow$  a random hyper-parameter combination
     $params\_mse \leftarrow 0$ 
    for  $k_i \leftarrow 1 \dots K_I$  do                             $\triangleright$  Evaluate  $params$  with  $K_I$ -fold CV
      Generate  $k$ -th inner fold  $(D_T^i, D_V^i)$  from  $D_T^o$ 
       $model \leftarrow$  a new instance of the model with parameters  $params$ 
      Train  $model$  on  $D_T^i$ 
       $mse \leftarrow$  result of the evaluation of  $model$  on  $D_V^i$ 
       $params\_mse \leftarrow params\_mse + mse$ 
    end for
    if  $params\_mse < best\_mse$  then                             $\triangleright$  MSE comparison on  $K$ -fold CV result
       $best\_mse \leftarrow params\_mse$ 
       $best\_params \leftarrow params$ 
    end if
  end for
   $model \leftarrow$  a new instance of the model with parameters  $best\_params$ 
  Train  $model$  on  $D_T^o$ 
  Evaluate  $model$  on  $D_V^o$  and compute evaluation metrics
end for
Compute summary of the metrics obtained in the outer loop

```

Chapter 4

Results

todo

4.1 Exploratory Data Analysis

The dataset consists of 17 years of measurements, from January 2001 to December 2017, for a total of 3 436 416 observations. Unfortunately, most of these observations contain low quality measurements, as detailed in section 3.1.5, leaving only 1 561 973 usable records. Additionally, the turbulent fluxes measured in March 2016 exhibit a much wider range than the March measurements of other years. Roughly 15 to 20% of the measurements in that month are suspicious; since the Cesar database contains a separate dataset every month, we decided to completely exclude the dataset of March 2016, fearing for a systematic error somewhere in the process. Finally, 6 more measurements of the turbulent latent heat are way above the acceptable range. This leaves 1 552 414 usable records.

4.2 Gradient Computation

compare methods for computing gradient

4.3 Monin-Obukhov Similarity Theory

performance of MOST in table 4.1

4.4 Model Comparison

performance of the models

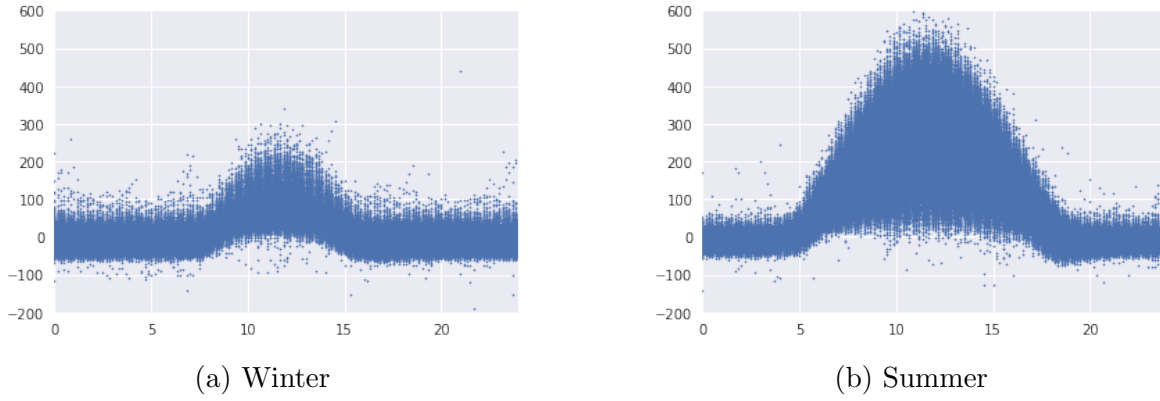


Figure 4.1: Turbulent latent heat flux versus hour of day in Winter (left) and Summer (right). One can see both the difference in day duration, and the effect of increased temperature on evaporation. The sensible heat flux follows a very similar pattern, with lower absolute values.

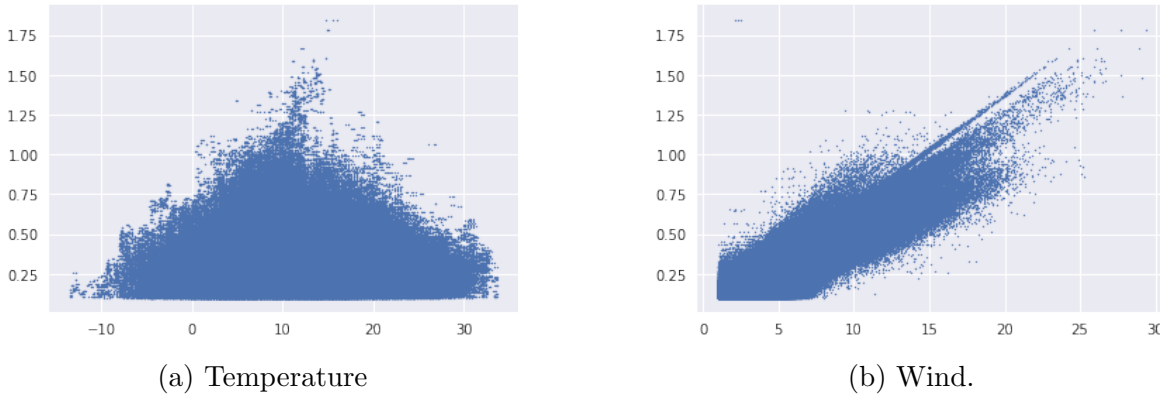


Figure 4.2: Friction velocity versus temperature (left) and wind speed (right). There is a clear subgroup where the wind and u^* are almost perfectly correlated; this is because of the gap-filling technique discussed in section 3.1.4. Additionally, there are three subgroups with different slopes of the regression line; they differ widely in fluxes (both surface and turbulent) and radiation.

Table 4.1: Evaluation metrics for the Monin-Obukhov estimator

	Mean	Std.
Mean Squared Error	0.672028	0.031320
R^2 Score	0.705471	0.008943
Mean Absolute Error	0.540632	0.013878
Median Absolute Error	0.371569	0.015601
Mean Absolute Percent Error	252.925389	196.086056
Median Absolute Percent Error	28.964847	0.938811

Table 4.2: Evaluation metrics for the Ridge linear regression estimator on F5 with trend

	Mean	Std.
Mean Squared Error	0.593038	0.048443
R^2 Score	0.734663	0.017526
Mean Absolute Error	0.509204	0.018504
Median Absolute Error	0.345660	0.009037
Mean Absolute Percent Error	264.613746	298.463690
Median Absolute Percent Error	27.816959	1.286721

Table 4.3: Evaluation metrics for the k-nearest neighbors estimator on F5 with trend

	Mean	Std.
Mean Squared Error	0.583102	0.035314
R^2 Score	0.739219	0.010517
Mean Absolute Error	0.487180	0.016797
Median Absolute Error	0.300725	0.011178
Mean Absolute Percent Error	151.571688	55.287673
Median Absolute Percent Error	25.450125	1.270630

Table 4.4: Evaluation metrics for the gradient boosting estimator on F5 with trend

	Mean	Std.
Mean Squared Error	0.285956	0.012074
R^2 Score	0.872323	0.006331
Mean Absolute Error	0.319557	0.006457
Median Absolute Error	0.180583	0.003956
Mean Absolute Percent Error	88.093397	53.138612
Median Absolute Percent Error	15.888091	0.499692

Chapter 5

Discussion

actual constants and coefficients (e.g. von karmans constant) dont really matter, because they are constant and do not affect the ml model

5.1 Limitations and Future Work

one of the reasons why MOST is unreliable in stable conditions is because of the difficulty of measuring fluxes

our models can only be used in conditions that are similar to cabauw. most importantly, they are not valid in the oceanic surface layer, because it is very different [citation]

Bibliography

- Sylvain Arlot, Alain Celisse, et al. A survey of cross-validation procedures for model selection. *Statistics surveys*, 4:40–79, 2010.
- Peter Baas, Bas Van de Wiel, Steven van der Linden, and F C. Bosveld. From near-neutral to strongly stratified: Adequately modelling the clear-sky nocturnal boundary layer at cabauw. 166, 10 2017.
- Yoshua Bengio and Yves Grandvalet. No unbiased estimator of the variance of k-fold cross-validation. In S. Thrun, L. K. Saul, and B. Schölkopf, editors, *Advances in Neural Information Processing Systems 16*, pages 513–520. MIT Press, 2004.
- James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *J. Mach. Learn. Res.*, 13(1):281–305, February 2012. ISSN 1532-4435. URL <http://dl.acm.org/citation.cfm?id=2503308.2188395>.
- L. Breiman, J. Friedman, R. Olshen, and C. Stone. *Classification and Regression Trees*. Wadsworth and Brooks, Monterey, CA, 1984.
- Leo Breiman. Random forests. *Machine Learning*, 45(1):5–32, Oct 2001. ISSN 1573-0565. doi: 10.1023/A:1010933404324. URL <https://doi.org/10.1023/A:1010933404324>.
- Hugh Chipman, Edward I. George, and Robert E. McCulloch. *The Practical Implementation of Bayesian Model Selection*, volume Volume 38 of *Lecture Notes–Monograph Series*, pages 65–116. Institute of Mathematical Statistics, Beachwood, OH, 2001. doi: 10.1214/lnms/1215540964. URL <https://doi.org/10.1214/lnms/1215540964>.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *J. Mach. Learn. Res.*, 12:2121–2159, July 2011. ISSN 1532-4435. URL <http://dl.acm.org/citation.cfm?id=1953048.2021068>.
- Bradley Efron and Robert J. Tibshirani. *An Introduction to the Bootstrap*. Number 57 in Monographs on Statistics and Applied Probability. Chapman & Hall/CRC, Boca Raton, Florida, USA, 1993.
- Seymour Geisser. The predictive sample reuse method with applications. *Journal of the American Statistical Association*, 70(350):320–328, 1975. doi: 10.1080/01621459.1975.10479865.

- I. Gitman, D. Dilipkumar, and B. Parr. Convergence Analysis of Gradient Descent Algorithms with Proportional Updates. *ArXiv e-prints*, January 2018.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *CoRR*, abs/1502.01852, 2015. URL <http://arxiv.org/abs/1502.01852>.
- Ulf Högström. Non-dimensional wind and temperature profiles in the atmospheric surface layer: A re-evaluation. *Boundary-Layer Meteorology*, 42(1):55–78, Jan 1988. ISSN 1573-1472. doi: 10.1007/BF00119875. URL <https://doi.org/10.1007/BF00119875>.
- Ulf Högström. Review of some basic characteristics of the atmospheric surface layer. *Boundary-Layer Meteorology*, 78(3):215–246, Mar 1996. ISSN 1573-1472. doi: 10.1007/BF00120937. URL <https://doi.org/10.1007/BF00120937>.
- AAM Holtslag JW Verkaik. Wind profiles, momentum fluxes and roughness lengths at cabauw revisited. 2006.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014. URL <http://arxiv.org/abs/1412.6980>.
- Ron Kohavi. A study of cross-validation and bootstrap for accuracy estimation and model selection. In *Proceedings of the 14th International Joint Conference on Artificial Intelligence - Volume 2*, IJCAI’95, pages 1137–1143, San Francisco, CA, USA, 1995. Morgan Kaufmann Publishers Inc. ISBN 1-55860-363-8. URL <http://dl.acm.org/citation.cfm?id=1643031.1643047>.
- Ann Korrell, H.A. Panosky, and R.J. Rossi. Wind profiles at the boulder tower. 08 1981.
- Xuhui Lee, William Massman, and Beverly Law. *Handbook of Micrometeorology: A Guide for Surface Flux Measurement and Analysis*. Kluwer Academic Publishers, 2004.
- W.A.A. Monna. Comparative investigation of dynamic properties of some propeller vanes. Technical report, Koninklijk Nederlands Meteorologisch Instituut, 1978.
- Arvind Neelakantan, Luke Vilnis, Quoc V. Le, Ilya Sutskever, Lukasz Kaiser, Karol Kurach, and James Martens. Adding gradient noise improves learning for very deep networks. *CoRR*, abs/1511.06807, 2015. URL <https://arxiv.org/abs/1511.06807>.
- Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer, New York, NY, USA, 1999.
- Michael Optis, Adam Monahan, and Fred C. Bosveld. Moving beyond monin-obukhov similarity theory in modelling wind-speed profiles in the lower atmospheric boundary layer under stable stratification. 153:497–514, 12 2014.

- Sebastian Ruder. An overview of gradient descent optimization algorithms. *CoRR*, abs/1609.04747, 2016. URL <http://arxiv.org/abs/1609.04747>.
- D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Parallel distributed processing: Explorations in the microstructure of cognition, vol. 1. chapter Learning Internal Representations by Error Propagation, pages 318–362. MIT Press, Cambridge, MA, USA, 1986. ISBN 0-262-68053-X. URL <http://dl.acm.org/citation.cfm?id=104279.104293>.
- R. S. Scorer. Atmospheric data analysis, roger daley, cambridge atmospheric and space science series, cambridge university press, cambridge, 1991. no. of pages: xiv + 457. isbn 0521 382157. *International Journal of Climatology*, 12(7):763–764, 1992. ISSN 1097-0088. doi: 10.1002/joc.3370120708. URL <http://dx.doi.org/10.1002/joc.3370120708>.
- Samuel L. Smith and Quoc V. Le. A bayesian perspective on generalization and stochastic gradient descent. *CoRR*, abs/1710.06451, 2017. URL <http://arxiv.org/abs/1710.06451>.
- C.J Stone. Consistent nonparametric regression. *The Annals of Statistics*, 5(4):595–620, July 1977.
- Mervyn Stone. Cross-validatory choice and assessment of statistical predictions. *Journal of the royal statistical society. Series B (Methodological)*, pages 111–147, 1974.
- F T. M. Nieuwstadt. The turbulent structure of the stable, nocturnal boundary layer. 41:2202–2216, 07 1984.
- N. K. Vinnichenko. The kinetic energy spectrum in the free atmosphere 1 second to 5 years. *Tellus*, 22(2):158–166, 1970. ISSN 2153-3490. doi: 10.1111/j.2153-3490.1970.tb01517.x. URL <http://dx.doi.org/10.1111/j.2153-3490.1970.tb01517.x>.
- Larry Wasserman. Bayesian model selection and model averaging. *Journal of Mathematical Psychology*, 44(1):92 – 107, 2000. ISSN 0022-2496. doi: <https://doi.org/10.1006/jmps.1999.1278>. URL <http://www.sciencedirect.com/science/article/pii/S0022249699912786>.
- H.R.A. Wessels. Distortion of the wind field by the cabauw meteorological tower. Technical report, Koninklijk Nederlands Meteorologisch Instituut, 1983.
- Matthew D. Zeiler. Adadelata: An adaptive learning rate method. *CoRR*, abs/1212.5701, 2012.
- Yongli Zhang and Yuhong Yang. Cross-validation for selecting a model selection procedure. *Journal of Econometrics*, 187(1):95–112, 7 2015. ISSN 0304-4076. doi: 10.1016/j.jeconom.2015.02.006.