ECS171: Machine Learning

L3 Validation of hypotheses: regression

Instructor: Prof. Maike Sonnewald TAs: Pu Sun & Devashree Kataria



Intended Learning Outcomes

- Describe bias and variance, as well as the bias-variance trade-off
- For regression models: Describe difference in the expression for quantifying errors and calculate these for data
- Describe mathematically and apply, for linear regression and multivariate linear regression the analytical (Ordinary Least Squares) and numerical (Gradient Descent) methods for determining fit

Rec. reading: B 3.1 + R 1-7, R 3.4.6

Resource: matrix_OLS_NYU_notes.pdf

Hypothesis testing: Goodness of fit

- Indicate the fit of a trained model and the causes of poor performance in machine learning.
- Goal of training a model: **Obtaining a** generalized model to perform well on unseen data.
- Overfitting:
 - A function (ML model) is too closely aligned with a limited set of data
- Underfitting:
 - A function (ML model) is not aligned with a limited dataset
- Model complexity:
 - As simple as possible, but no simpler



A model is a function that represents the data Model 100 0.4 0.2 0 -0.2 -0.4 -0.6 -0.5 -0.4 -0.3 -0.2 -0.1 0 0.1 0.2 20 30 40 50 60 70 80 90 100 Classification Regression

Hypothesis: If I fit an ML model it will mimic the underlying 'model' the data came from

Predicting continuous outputs: Regression

We need:

- Features (inputs): we'll call these x (or **x** if vectors)
- Training examples: many xⁱ for which yⁱ is known
- A **model:** Function that represents the relationship between x and y

$$y(x) =$$
function (x, \mathbf{w})

Linear: $y(x) = w_0 + w_1$

- A **loss function**: Estimate how well our model approximates the training examples (aka cost or objective function)
- **Optimization**, a way of finding the parameters of our model that minimizes the loss function



The difference between the 'true' function and our model is estimated using our observations



Symptoms of over and underfitting

Underfitting:

- High training error
- Training and test error similar
- High bias

Overfitting:

- Very low training error
- Test error much higher training error
- High variance on test data



How do we measure a model's performance?

1. Regression Metrics (Continuous values)

- Sum of Squares Error (SSE): Measures the total deviation of the response values from the fit to the response values.
- Mean Absolute Error (MAE): The average of the absolute differences between the predicted values and actual values.
- Mean Squared Error (MSE): The average of the squared differences between the predicted values and actual values.
- **Root Mean Squared Error (RMSE)**: The square root of MSE, often more interpretable as it is in the same units as the response variable.
- **R-squared (Coefficient of Determination)**: Measures the proportion of the variance in the dependent variable that is predictable from the independent variables.

2. Classification Metrics (Categorize data into labels)

- Accuracy: The proportion of total predictions that were correct.
- **Precision**: The proportion of positive identifications that were actually correct.
- Recall (Sensitivity): The proportion of actual positives that were identified correctly.
- **F1 Score**: The harmonic mean of precision and recall.
- **Confusion Matrix**: A table used to describe the performance of a classification model, showing the actual vs. predicted values.
- **ROC-AUC**: The area under the receiver operating characteristic curve, measuring the trade-off between true positive rate and false positive rate.
- **Precision-Recall Curve**: Focuses on the performance with respect to the positive (minority) class.

3. Clustering Metrics (Grouping data into clusters)

- **Silhouette Score**: Measures how similar an object is to its own cluster compared to other clusters.
- **Davies-Bouldin Index**: The average 'similarity' between each cluster and its most similar cluster, where lower values indicate better clustering.
- **Calinski-Harabasz Index**: The ratio of the sum of between-clusters dispersion and of within-cluster dispersion.

4. Time Series Metrics

- **Mean Absolute Percentage Error (MAPE)**: The mean absolute percentage difference between the predicted and actual values.
- Symmetric Mean Absolute Percentage Error (sMAPE): An adjustment to MAPE that handles zero values more effectively.

5. Other Methods

- **Cross-Validation**: A method to assess the robustness of the model by training and testing the model on different subsets of the dataset.
- Learning Curves: Plotting the model performance on the training set and the validation set over time or over the number of datasets.
- **Feature Importance**: Evaluating which features contribute most to the model's predictive power.

Too simple: inflexible learning due to too few/wrong features or too strict regularization \rightarrow little variance but more bias

Too complex: more prediction variance



Quantifying errors:

- \hat{y}_i is the predicted value for the i^{th} observation.
- y_i is the true value for the i^{th} observation.
- $\bullet n$ is the total number of observations.
- * $ar{y}$ is the mean of all observed values of the dependent variable.

Bias is the average of the errors predictions made by the model and the true values

$$ext{Bias} = rac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)$$

Sum of Squares Error: Squaring avoids cancellation of pos and neg and emphasizes larger errors. Aka Residual Sum of Squares Error (RSS)

$$ext{SSE} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Mean Squared Error: Averaged SSE, giving an average error per data point

$$ext{MSE} = rac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Quantifying errors cont.

- \hat{y}_i is the predicted value for the i^{th} observation.
- y_i is the true value for the i^{th} observation.
- *n* is the total number of observations.
- * $ar{y}$ is the mean of all observed values of the dependent variable.

Mean Absolute Error: Avoids squaring errors, useful if large errors are not worse than smaller ones

$$ext{MAE} = rac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Root Mean Squared Error: Same units as response variable, good when large errors are problems

$$ext{RMSE} = \sqrt{rac{1}{n}\sum_{i=1}^n(y_i-\hat{y}_i)^2}$$

Coefficient of Determination (R²): Measures of how well the independent variables explain the variability in the dependent variable in a regression model.

$$R^2 = 1 - rac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - ar{y})^2}$$

Bias and variance

Bias

- **Definition**: Bias refers to the error due to overly simplistic assumptions in the learning algorithm. It can lead to the model underfitting the training data, meaning it does not capture the underlying trends well.
- **Characteristics**: A high-bias model is likely to have a lower level of complexity, which makes it less flexible in learning from the data. This results in the model missing relevant relations between features and target outputs.
- **Consequence**: High bias can cause the model to be less accurate on both training and testing data, leading to poor generalization.



Bias and variance

Variance

- **Definition**: Variance refers to the error due to too much complexity in the learning algorithm. It can lead to the model overfitting the training data, meaning it captures noise along with the underlying patterns.
- **Characteristics**: A high-variance model is extremely sensitive to the fluctuations in the training data. It learns features from the training data that don't generalize to unseen data.
- **Consequence**: High variance can cause the model to perform well on the training data but poorly on unseen (test) data due to overfitting.

Variance is measured through, for example, the difference between SSE for the training and the testing data. If $SSE_{test} - SSE_{train}$ is high the model has high variance



Bias-variance trade-off

- **Bias** is about the simplicity of the model high bias can lead to underfitting.
- Variance is about the complexity of the model high variance can lead to overfitting.
- Effective machine learning involves managing this tradeoff to achieve a model that generalizes well to new, unseen data.



Linear Regression (LR) fits a linear function

Goal to find function f so that:

f(x) = y

Two approaches to determine f(x):

- Analytical: Ordinary Least Squares (OLS)
- Numerical: Gradient Descent



Linear Regression (LR) fits a linear function

One variable Consider a dataset with *n* observations (x_i, y_i) , where y_i is the $y_i = \beta_0 + \beta_1 x_i + \varepsilon$ dependent variable and x_i is the independent variable. The linear Observed Value of Y for X regression model is: em Error Predicted Value of $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ Y for X Slope = / We aim to find the values of β_0 (intercept) and β_1 (slope) that minimize the sum of squared residuals (errors), where the residual for each Intercept = β_0 observation is $\epsilon_i = (y_i) - (\beta_0 + \beta_1 x_i)$. X The 'true' value is: $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ Bias The predicted value is: $\widehat{y}_i = eta_0 + eta_1 x_i$

Estimating goodness-of-fit

- The fit of the model is estimated looking at the errors



Best line: - bias: -5.25135e-16 - SSE: 99.2439 - MSE: 0.992439 - MAE: 0.849258 - RMSE: 0.996212 - R²: 0.935972

Example: example_linear_regression.ipynb

Multivariate Linear Regression

Basic equation for a multivariate linear regression model:

 $y_i=w_0+w_1x_{i1}+w_2x_{i2}+\dots+w_px_{ip}+\epsilon_i$

 y_i is the dependent variable for the i^{th} observation. $x_{i1}, x_{i2}, \ldots, x_{ip}$ are the independent variables (predictors) for the i^{th} observation. $w_0, w_1, w_2, \ldots, w_p$ are the coefficients to be estimated. ϵ_i is the error term for the i^{th} observation.

The w_o is the bias aka the intercept

Model achieves a good fit for the regression line by finding the best coefficient values (w) that **minimize the errors**.

Indepevariabl	Target Variable (\	Target Variable (Y)		
Temperature		Yield		
		142		
		149		
		161		
72		167		
		168		
	52	162		
76		171		
79	52	175		
80	62	182		

A $\in \mathbb{R}^{(73,2.5)}$ With many variables (predictors) we fit (hyper)planes

Y

Each predictor (i) has an associates slope (w)

For each dependent variable (data point i), with observations indexed by p: $\underline{x}_i = x_{i1} + x_{i2} + \dots + x_{ip}$

We incorporate the bias w_0 by setting $x_0=1$, so the prediction for n observations becomes:

$$y(\mathbf{x}) = w_0 + \sum_{i=1}^n w_i x_i = \sum_{i=0}^n w_i \mathbf{x}_i$$

The true value includes an error term of the residuals: $\hat{y}(\mathbf{x}) = w_0 + \sum_{i=1}^n w_i x_i + \epsilon_i$ Using matrix notation: $\bigvee \in \mathbb{R}^{n \times (p+1)} \times \mathbb{E} = \mathbb{R}^{n \times 1} \times \mathbb{E} = \mathbb{E} = \mathbb{R}^{n \times 1} \times \mathbb{E} = \mathbb{E}$



The cost function is used to minimize the error

- The cost function (J) is here defined to minimize the sum of squared errors (SSE) by varying the coefficients w:

$$J(\mathbf{w}) = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

- Recall, the predicted value is: $\hat{y_i} = w_0 + w_1 x_{i1} + w_2 x_{i2} + \dots + w_p x_{ip}$
- By minimizing J(**w**) the relationship between y_i and \hat{y}_i can be approximated in the best available way
- Does not have to be SSE

(1)

Fitting the model analytically: Ordinary Least Squares (OSL)

- Differentiation finds the point at which a function f(x) is at a minimum denoted f'(x)
- We can differentiate the cost function $J(\mathbf{w})$ wrt **w** to find our minimum: $\frac{\partial J(w)^{\ell}}{\partial w} = 0$



Matrix notation

Given a dataset with \boldsymbol{n} observations and \boldsymbol{p} independent variables, the model is:

 $\mathbf{Y} = \mathbf{X}\mathbf{W} + \epsilon$

Where:

- * ${f Y}$ is an n imes 1 column vector of the dependent variable.
- * ${f X}$ is an n imes (p+1) matrix of the independent variables, with the first column as 1s for the intercept term.
- * ${f W}$ is a (p+1) imes 1 column vector of coefficients (including the intercept).
- * ϵ is an $n\times 1$ column vector of the residuals.

The response vector ${f Y}$ is:

$$\mathbf{Y} = egin{bmatrix} y_1 \ y_2 \ dots \ y_n \end{bmatrix}$$

• y_i represents the value of the dependent variable for the i^{th} observation.

The data matrix ${f X}$ is structured as follows:

	$\lceil 1 \rceil$	x_{11}	x_{12}		x_{1p}
V	1	x_{21}	x_{22}		x_{2p}
$\mathbf{X} =$:	:	÷	·	:
	1	x_{n1}	x_{n2}		x_{np}

- Each row represents an observation.
- The first column is all 1's for the intercept term.
- * x_{ij} represents the value of the j^{th} predictor for the i^{th} observation.
- * n is the number of observations, and p is the number of predictors.

The coefficient vector ${\boldsymbol{W}}$ is structured as:



* w_0 is the intercept term.

• w_1,\ldots,w_p are the coefficients of the predictors.

Minimize a cost function to fit the model finding optimal w



Fitting the model analytically: Ordinary Least Squares (OSL) FOIL 322 In matrix notation (J(w) is now SSE): $\nabla_{\mathbf{W}} \mathscr{B}_{\mathbf{K}} = \nabla_{\mathbf{W}} [\mathbf{Y}^T \mathbf{Y} - \mathbf{Y}^T \mathbf{X} \mathbf{W} - \mathbf{W}^T \mathbf{X}^T \mathbf{Y} + \mathbf{W}^T \mathbf{X}^T \mathbf{X} \mathbf{W}]$ Taking the derivative term by term, and noting that $Y^{T}Y$ is a constant with respect to W, and - $\overline{C} \overline{V}_{w} \left[- \overline{V}_{x}^{T} Y - \overline{V}_{x}^{T} Y + \overline{V}_{x}^{T} Y \right]$ $\overline{C} \overline{V}_{w} \left[- 2 \overline{V}_{x}^{T} Y + \overline{V}_{x}^{T} Y \right]$ $Y^{T}XW$ and $W^{T}X^{T}Y$ are equivalent, the derivative simplifies to: (4) $\nabla_{\mathbf{W}}$ size $= -2\mathbf{X}^T\mathbf{Y} + 2\mathbf{X}^T\mathbf{X}\mathbf{W}$ $-\mathbf{\mathbf{Z}}\mathbf{X}^{T}\mathbf{Y} + \mathbf{\mathbf{Z}}\mathbf{X}^{T}\mathbf{X}\mathbf{W} = 0$ Set to zero: (KTXW = XTY Rearranging: Finally:

- This result gives you the OLS estimator for the coefficients in multiple linear regression.

Fitting the model Numerically: Gradient Descent

- If we can't find an analytical solution we use a numerical method to find optimal weights **w**
- Gradient descent is iterative (do many times with updates):
 - Initialise w.e.g. with random numbers
 - Calculate J(w), change w_i, ask: Has the error J(w) gotten smaller?
 - We update w_p : Is it moving towards an optimum?

It can also be a 'hill-climber'. Take note of the sign of the cost function





Fitting the model Numerically: Gradient Descent

- How do we know which direction to move the weights?
 - Assess the gradient resulting from changing the weights for one observation (w_i)





Least Mean Squares (LMS) update rule



- As error approaches zero, so does the update (w changes less)

Optimizing across a training set

- To generalise to all data points in the training set options include:

- Full 1: Batch updates: Sum or average updates across every example i, then change the parameter values: Batch $\mathbf{w} \leftarrow \mathbf{w} + 2\lambda \sum_{i=1}^{n} (y_i - \hat{y}_i(x_i)) \underline{x}_i$ Mini-batch $\mathbf{w} \leftarrow \mathbf{w} + 2\lambda \sum_{i=1}^{n} (y_i - \hat{y}_i(x_i)) \underline{x}_i$
 - 2: <u>Stochastic/online</u> updates: Update the parameters for each data point in turn, according to its own gradients

 Algorithm 1 Stochastic gradient descent

 1: Randomly shuffle examples in the training set

 2: for i = 1 to \mathbf{v} do

 3: Update:

 $\mathbf{w} \leftarrow \mathbf{w} + 2\lambda(y_i - \hat{y}_i(x_i))x_i$ (update for a linear model)

 4: end for

Gradient Descent notes

- The learning rate (λ) represents the size of the 'steps' of the descent
- Pros:
 - Intuitive
 - Heuristic approach (stochastic optimisation)

Con ex

- Cons:
 - Can take many iterations to converge
 - Only optimal for <u>'smooth</u>' functions

We will get back to working through an example of the GD Update Rule for 1 observation



The choice of loss function can significantly impact results



Example: example_linear_regression.ipynb