Recap: regression Statistical Methods in NLP 2 ISCL-BA-08

#### Çağrı Çöltekin ccoltekin@sfs.uni-tuebingen.de

University of Tübingen Seminar für Sprachwissenschaft

Summer Semester 2025

version: 17f8995 @2025-04-30

# Regression

- In ML, regression refers to a problem where outcome variable is a number (numeric quantity, continuous random variable)
- In practice, regression most commonly refers to *linear regression*, solved using *least squares* optimization
- We will review solving the linear regression problem from different viewpoints

Linear regression is about finding a linear *model* of the form,

$$\mathbf{y} = w_1 \mathbf{x} + w_0$$

- y is a numeric quantity we want to predict
- x is a measurement/value helpful for predicting y
- *w*<sup>0</sup> and *w*<sup>1</sup> are the parameters that we want to learn from data
- both x and y can be vector valued



Linear regression is about finding a linear *model* of the form,

$$y = w_1 x + w_0$$

- y is a numeric quantity we want to predict
- x is a measurement/value helpful for predicting y
- *w*<sup>0</sup> and *w*<sup>1</sup> are the parameters that we want to learn from data
- both x and y can be vector valued



Linear regression is about finding a linear *model* of the form,

$$y = w_1 x + w_0$$

- y is a numeric quantity we want to predict
- x is a measurement/value helpful for predicting y
- *w*<sup>0</sup> and *w*<sup>1</sup> are the parameters that we want to learn from data
- both x and y can be vector valued



Linear regression is about finding a linear *model* of the form,

$$\mathbf{y} = w_1 \mathbf{x} + w_0$$

- y is a numeric quantity we want to predict
- x is a measurement/value helpful for predicting y
- *w*<sup>0</sup> and *w*<sup>1</sup> are the parameters that we want to learn from data
- both x and y can be vector valued



### Estimating regression parameters

- We view learning as a search for the regression equation with least error
- The error terms are also called *residuals*
- We want error to be low for the whole training set: average (or sum) of the error has to be reduced
- Can we minimize the sum of the errors?



### Least squares regression

In least squares regression, we want to find  $w_0$  and  $w_1$  values that minimize

$$E(\boldsymbol{w}) = \sum_{i} (y_i - (w_0 + w_1 x_i))^2$$

- Note that E(w) is a *quadratic* function of  $w = (w_0, w_1)$
- As a result, E(w) is *convex* and have a single extreme value
  - there is a unique solution for our minimization problem
- In case of least squares regression, there is an analytic solution
- Even if we do not have an analytic solution, if the error function is convex, a search procedure like *gradient descent* can still find the *global minimum*

# A simple example

• Data: 
$$\mathbf{x} = \begin{bmatrix} 4 \\ 2 \end{bmatrix}$$
  $\mathbf{y} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$   
Model:  $\widehat{\mathbf{y}} = w\mathbf{x}$ 

# A simple example

• Data: 
$$\mathbf{x} = \begin{bmatrix} 4 \\ 2 \end{bmatrix}$$
  $\mathbf{y} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$   
Model:  $\hat{\mathbf{y}} = w\mathbf{x}$ 

• Squared errors

$$E(w) = (4w - 1)^{2} + (2w - 2)^{2}$$
$$= 20w^{2} - 16w + 5$$



# A simple example

• Data: 
$$\mathbf{x} = \begin{bmatrix} 4 \\ 2 \end{bmatrix}$$
  $\mathbf{y} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$   
Model:  $\hat{\mathbf{y}} = w\mathbf{x}$ 

• Squared errors

$$E(w) = (4w - 1)^{2} + (2w - 2)^{2}$$
$$= 20w^{2} - 16w + 5$$

• Setting the derivative to zero:

$$\frac{\mathrm{dE}}{\mathrm{d}w} = 40w - 16 = 0 \Rightarrow w = \frac{2}{5}$$



# Regression with multiple variables

- The example generalizes to more parameters
- Instead of derivative with respect to a single variable, we calculate the gradient
- The solution is where the gradient is  $\boldsymbol{0}$
- This leads to a system of linear equations, whose solution vector is the best parameters

# Maximum Likelihood Estimation (MLE)

- In MLE the task is to find the model m that assigns the maximum *likelihood* to the observed data x
- To emphasize that likelihood is a function of model parameters, w, we indicate it as  $\mathcal{L}(w;x)$
- Formally, the task is finding

$$\boldsymbol{w}_{\mathrm{MLE}} = \operatorname*{arg\,max}_{\boldsymbol{w}} \mathcal{L}(\boldsymbol{w}; \boldsymbol{x})$$

• In most cases, working with log likelihood is easier, since log is a monotonically increasing function,

$$\boldsymbol{w}_{\text{MLE}} = \operatorname*{arg\,max}_{\boldsymbol{w}} \log \mathcal{L}(\boldsymbol{w}; \boldsymbol{x}) = \operatorname*{arg\,min}_{\boldsymbol{w}} - \log \mathcal{L}(\boldsymbol{w}; \boldsymbol{x})$$

# MLE for simple regression

 $\mathbf{y}_{\mathbf{i}} = \mathbf{w}_0 + \mathbf{w}_1 \mathbf{x}_{\mathbf{i}} + \mathbf{\varepsilon}_{\mathbf{i}}$ 

where  $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma)$ 

- We additionally assume that  $\boldsymbol{\sigma}$  is independent of  $\boldsymbol{x}$
- This means  $y \sim \mathcal{N}(w_0 + w_1 x, \sigma)$
- Now the likelihood function becomes,





MLE for simple regression (2)

Log likelihood: 
$$-n \ln \sigma \sqrt{2\pi} - \frac{1}{2\sigma^2} \sum_{i=1}^{n} (y_i - (w_0 + w_1 x_i))^2$$

• Note that maximizing log likelihood is equivalent to minimizing

$$\sum_{i=1}^{n} (y_i - (w_0 + w_1 x_i))^2$$

- This is the squared error (the same as what we did before)
- MLE estimate of the regression parameters is equivalent to least-squares regression

#### Approximate solutions to systems of linear equations

$$\begin{bmatrix} 4 & 2 \\ 2 & 2 \\ 4 & 3 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 10 \\ 4 \\ 4 \end{bmatrix}$$

• Can we solve the equation above?

#### Approximate solutions to systems of linear equations

$$\begin{bmatrix} 4 & 2 \\ 2 & 2 \\ 4 & 3 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 10 \\ 4 \\ 4 \end{bmatrix}$$

- Can we solve the equation above?
- Can we find the 'best' approximation?

### Approximate solutions to systems of linear equations

$$\begin{bmatrix} 4 & 2 \\ 2 & 2 \\ 4 & 3 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 10 \\ 4 \\ 4 \end{bmatrix}$$

- Can we solve the equation above?
- Can we find the 'best' approximation?
- Reminder: finding the solution means

$$\begin{bmatrix} 4\\2\\4 \end{bmatrix} w_1 + \begin{bmatrix} 2\\2\\3 \end{bmatrix} w_2 = \begin{bmatrix} 10\\4\\4 \end{bmatrix}$$

# Picture of the (non)solution

- In higher dimensional spaces we want the projection onto the column space of X
- The error vector *e* is perpendicular to all column vectors of X, x<sub>i</sub>
- Again, note that e = y p



# Deriving linear regression on higher dimensions

$$\begin{split} & X^{\mathsf{T}}(\mathbf{y} - \mathbf{p}) = 0 & \text{Error vector is orthogonal to columns} \\ & X^{\mathsf{T}}(\mathbf{y} - \mathbf{X} \mathbf{w}) = 0 & \mathbf{p} \text{ is the weighted combination of columns} \\ & X^{\mathsf{T}} \mathbf{X} \mathbf{w} = \mathbf{X}^{\mathsf{T}} \mathbf{y} & \text{Note: } \mathbf{X}^{\mathsf{T}} \mathbf{X} \text{ is square} \\ & \mathbf{w} = (\mathbf{X}^{\mathsf{T}} \mathbf{X})^{-1} \mathbf{X}^{\mathsf{T}} \mathbf{y} & \text{The final solution} \end{split}$$

The projection of **y** onto columns space of **X** is

$$\mathbf{p} = \mathbf{X}(\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$$

Ç. Çöltekin, SfS / University of Tübingen

$$\times \begin{bmatrix} a & 0 & 0 & 0 \\ 0 & b & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} =$$

$$\begin{bmatrix} 1/a & 0 & 0\\ 0 & 1/b & 0\\ 0 & 0 & 0\\ 0 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} a & 0 & 0 & 0\\ 0 & b & 0 & 0\\ 0 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0\\ 0 & 1 & 0 & 0\\ 0 & 0 & 0 & 0\\ 0 & 0 & 0 & 0 \end{bmatrix}$$

• We want matrix multiplication to get as close to I as possible. Consider the  $3 \times 4$  diagonal matrix:

$$\begin{bmatrix} 1/a & 0 & 0\\ 0 & 1/b & 0\\ 0 & 0 & 0\\ 0 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} a & 0 & 0 & 0\\ 0 & b & 0 & 0\\ 0 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0\\ 0 & 1 & 0 & 0\\ 0 & 0 & 0 & 0\\ 0 & 0 & 0 & 0 \end{bmatrix}$$

• For an  $n \times n$  diagonal matrix  $\Sigma$ ,  $\Sigma^+ = \Sigma^{-1}$ 

$$\begin{bmatrix} 1/a & 0 & 0\\ 0 & 1/b & 0\\ 0 & 0 & 0\\ 0 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} a & 0 & 0 & 0\\ 0 & b & 0 & 0\\ 0 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0\\ 0 & 1 & 0 & 0\\ 0 & 0 & 0 & 0\\ 0 & 0 & 0 & 0 \end{bmatrix}$$

- For an  $n \times n$  diagonal matrix  $\Sigma$ ,  $\Sigma^+ = \Sigma^{-1}$
- For any invertible  $n \times n$  matrix **X**,  $X^+ = X^{-1}$

$$\begin{bmatrix} 1/a & 0 & 0\\ 0 & 1/b & 0\\ 0 & 0 & 0\\ 0 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} a & 0 & 0 & 0\\ 0 & b & 0 & 0\\ 0 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0\\ 0 & 1 & 0 & 0\\ 0 & 0 & 0 & 0\\ 0 & 0 & 0 & 0 \end{bmatrix}$$

- For an  $n \times n$  diagonal matrix  $\Sigma$ ,  $\Sigma^+ = \Sigma^{-1}$
- For any invertible  $n \times n$  matrix **X**,  $X^+ = X^{-1}$
- Singular value decomposition (SVD) provides the general solution:

$$X^+ = V \Sigma^+ U^T$$

# Regression through pseudo inverse

- Pseudo inverse is another method to find the regression parameters
- We want

$$Xw = y$$

but there is no general solution. Multiplying both sides with the pseudo inverse results in the best approximation

$$X^+ X w \approx X^+ y$$
$$w \approx X^+ y$$

• This also allows regression with 'wide' matrices, in which case we get lowest L2-norm solution

#### Final remarks

- Regression is probably the most popular method for all (scientific) research
- Many statistical methods are variations/extensions of regression
- Regression is also part of almost any ML method

### Final remarks

- Regression is probably the most popular method for all (scientific) research
- Many statistical methods are variations/extensions of regression
- Regression is also part of almost any ML method

Next:

• Recap: classification / evaluation

### Some sources of information

- Any modern linear algebra book (e.g., Strang, 2009) would cover regression
- For a more ML-focused introductions, also see James et al. (2024) or any machine learning textbook

- James, Gareth, Daniela Witten, Trevor Hastie, Robert Tibshirani, and Jonathan Taylor (2024). An introduction to statistical learning. Springer. ISBN: 9783031391897. URL: https://www.statlearning.com/.
- Strang, Gilbert (2009). Introduction to Linear Algebra, Fourth Edition. 4th ed. Wellesley Cambridge Press. ISBN: 9780980232714.