

Lecture 12 | 12.05.2025

Regression models beyond linearity

Linear regression models

□ Normal linear regression model

- generic regression model $Y = \mathbf{X}^\top \beta + \varepsilon$, for $\varepsilon \sim N(0, \sigma^2)$
- random sample $\{(Y_i, \mathbf{X}_i^\top)^\top; 1 = 1, \dots, n\}$ from $F_{(Y, \mathbf{X})}$
- conditional distribution of $Y|\mathbf{X}$ is normal, i.e., $Y|\mathbf{X} \sim N(\mathbf{X}^\top \beta, \sigma^2)$
- parameter estimates (LSE/MLE) are BLUE and normally distributed

□ Linear regression model without normality

- generic regression model $Y = \mathbf{X}^\top \beta + \varepsilon$, $\varepsilon \sim (0, \sigma^2)$, $E\varepsilon^2 = \sigma^2 \in (0, \infty)$
- mean ($E[Y|\mathbf{X}] = \mathbf{X}^\top \beta$) and variance ($\text{Var}(Y|\mathbf{X}) = \sigma^2(\mathbf{X})$) specification
- conditional distribution of $Y|\mathbf{X}$ is left unspecified (LSE only)
- parameter estimates (MLE) are BLUE and asymptotically normal

Linear regression models

□ Normal linear regression model

- generic regression model $Y = \mathbf{X}^\top \beta + \varepsilon$, for $\varepsilon \sim N(0, \sigma^2)$
- random sample $\{(Y_i, \mathbf{X}_i^\top)^\top; 1 = 1, \dots, n\}$ from $F_{(Y, \mathbf{X})}$
- conditional distribution of $Y|\mathbf{X}$ is normal, i.e., $Y|\mathbf{X} \sim N(\mathbf{X}^\top \beta, \sigma^2)$
- parameter estimates (LSE/MLE) are BLUE and normally distributed

□ Linear regression model without normality

- generic regression model $Y = \mathbf{X}^\top \beta + \varepsilon$, $\varepsilon \sim (0, \sigma^2)$, $E\varepsilon^2 = \sigma^2 \in (0, \infty)$
- mean ($E[Y|\mathbf{X}] = \mathbf{X}^\top \beta$) and variance ($\text{Var}(Y|\mathbf{X}) = \sigma^2(\mathbf{X})$) specification
- conditional distribution of $Y|\mathbf{X}$ is left unspecified (LSE only)
- parameter estimates (MLE) are BLUE and asymptotically normal

Recall, that **linearity** + **normality** = "**lightness of being**" but linear regression models without the assumptions of normality introduce just a minor complication...

Linear regression models

□ Normal linear regression model

- generic regression model $Y = \mathbf{X}^\top \beta + \varepsilon$, for $\varepsilon \sim N(0, \sigma^2)$
- random sample $\{(Y_i, \mathbf{X}_i^\top)^\top; 1 = 1, \dots, n\}$ from $F_{(Y, \mathbf{X})}$
- conditional distribution of $Y|\mathbf{X}$ is normal, i.e., $Y|\mathbf{X} \sim N(\mathbf{X}^\top \beta, \sigma^2)$
- parameter estimates (LSE/MLE) are BLUE and normally distributed

□ Linear regression model without normality

- generic regression model $Y = \mathbf{X}^\top \beta + \varepsilon$, $\varepsilon \sim (0, \sigma^2)$, $E\varepsilon^2 = \sigma^2 \in (0, \infty)$
- mean ($E[Y|\mathbf{X}] = \mathbf{X}^\top \beta$) and variance ($\text{Var}(Y|\mathbf{X}) = \sigma^2(\mathbf{X})$) specification
- conditional distribution of $Y|\mathbf{X}$ is left unspecified (LSE only)
- parameter estimates (MLE) are BLUE and asymptotically normal

Recall, that **linearity + normality = "lightness of being"** but linear regression models without the assumptions of normality introduce just a minor complication...

Thus, the linearity property is way more crucial!

(linearity of the predictor, linearity of the least squares, linearity of the expectation, linearity of the normal distribution, ...)

Beyond linearity

- ❑ In practice, however: **The truth is (almost) never linear!**
(however, the linearity assumption is a good and easy approximation)
- ❑ **What to do, when the linearity assumption fails?**
(the answer usually depends on the reason why the linearity fails)
- ❑ **Note that there are a few levels of linearity in the model**
(linearity of the predictor, linearity of the expectation, linearity of LS, ...)

Beyond linearity

- ❑ In practice, however: **The truth is (almost) never linear!**
(however, the linearity assumption is a good and easy approximation)
- ❑ **What to do, when the linearity assumption fails?**
(the answer usually depends on the reason why the linearity fails)
- ❑ **Note that there are a few levels of linearity in the model**
(linearity of the predictor, linearity of the expectation, linearity of LS, ...)
 - ❑ the data are too flexible (higher order approximations/splines)
 - ❑ the data are too irregular (piecewise approximation)
 - ❑ the data are too complex (additive models)
 - ❑ the data are too volatile (robust estimation approaches)
 - ❑ the nature of Y contradicts the linear model (GLM)
 - ❑ and many more reasons (and way more alternatives)

two generalizations beyond linearity

□ Linearity of the predictor

- linear predictor in terms of $\mathbf{X}^\top \beta$, where $\beta \in \mathbb{R}^p$ are unknown parameters
- the linear predictor is directly associated with the (theoretical) quantity of interest – the (conditional) expectation of Y (i.e., $E[Y|\mathbf{X}] = \mathbf{X}^\top \beta$)
- however, this direct association may not be realistic in some situations
- \implies **generalized linear models (GLM)** & **non-linear models (NLS)**

two generalizations beyond linearity

□ Linearity of the predictor

- linear predictor in terms of $\mathbf{X}^\top \beta$, where $\beta \in \mathbb{R}^p$ are unknown parameters
- the linear predictor is directly associated with the (theoretical) quantity of interest – the (conditional) expectation of Y (i.e., $E[Y|\mathbf{X}] = \mathbf{X}^\top \beta$)
- however, this direct association may not be realistic in some situations
- \implies **generalized linear models (GLM) & non-linear models (NLS)**

□ Linearity of the expectation

- the expectation $EY = \int_{\mathbb{R}} x dF_Y(x)$ of some random variable $Y \sim F_Y$ is a linear functional
- the expectation is also one of the most important characteristics of some unknown population (random variable)
- on the other hand, the expectation offers only a very limited information about the behavior of $Y \sim F_Y$
- \implies **quantile regression, expectile regression, or m-regression in general**

1. Generalized linear models

So far, all regression models concerned the response variable $Y \in \mathbb{R}$ that was a priori assumed to be continuous and the conditional distribution of $Y|\mathbf{X}$ was assumed to be normal or, at least, close to normal...

In practical applications, however, the domain of Y can be also more restricted...

- $Y \in \mathbb{N} \cup \{0\}$ (counts)
- $Y \in \{1, \dots, K\}$ for $K \in \mathbb{N}$ (categories/label)
- $Y \in \{0, 1\}$ (true/false)
- ...

1. Generalized linear models

So far, all regression models concerned the response variable $Y \in \mathbb{R}$ that was a priori assumed to be continuous and the conditional distribution of $Y|\mathbf{X}$ was assumed to be normal or, at least, close to normal...

In practical applications, however, the domain of Y can be also more restricted...

- $Y \in \mathbb{N} \cup \{0\}$ (counts)
- $Y \in \{1, \dots, K\}$ for $K \in \mathbb{N}$ (categories/label)
- $Y \in \{0, 1\}$ (true/false)
- ...

Note, that despite the fact that the domain of Y is restricted, the mean parameter of Y (the conditional mean if $Y|\mathbf{X}$ respectively) is still assumed to be from some compact subset, $\mathcal{M} \subset \mathbb{R}$... This is very useful in the following models...

Linear models with a flavour of nonlinearity

- in a **standard linear model**, the conditional mean is modelled as

$$E[Y|\mathbf{X}] = \mathbf{X}^T \boldsymbol{\beta}, \quad \text{for } \boldsymbol{\beta} \in \mathbb{R}^p$$

while the variance structure $\text{Var}[Y|\mathbf{X}]$ is modeled separately from the mean structure (e.g., $\text{Var}[Y|\mathbf{X}] = \sigma^2 \mathbb{I}$)

Linear models with a flavour of nonlinearity

- in a **standard linear model**, the conditional mean is modelled as

$$E[Y|\mathbf{X}] = \mathbf{X}^\top \boldsymbol{\beta}, \quad \text{for } \boldsymbol{\beta} \in \mathbb{R}^p$$

while the variance structure $\text{Var}[Y|\mathbf{X}]$ is modeled separately from the mean structure (e.g., $\text{Var}[Y|\mathbf{X}] = \sigma^2 \mathbb{I}$)

- in a **generalized linear model**, the conditional mean is modelled as

$$g(E[Y|\mathbf{X}]) = \mathbf{X}^\top \boldsymbol{\beta}, \quad \text{for } \boldsymbol{\beta} \in \mathbb{R}^p$$

for some **non-linear link function** $g : \mathcal{M} \rightarrow \mathbb{R}$ (typically continuous, smooth, regular, but nonlinear)

Linear models with a flavour of nonlinearity

- in a **standard linear model**, the conditional mean is modelled as

$$E[Y|\mathbf{X}] = \mathbf{X}^\top \boldsymbol{\beta}, \quad \text{for } \boldsymbol{\beta} \in \mathbb{R}^p$$

while the variance structure $\text{Var}[Y|\mathbf{X}]$ is modeled separately from the mean structure (e.g., $\text{Var}[Y|\mathbf{X}] = \sigma^2 \mathbb{I}$)

- in a **generalized linear model**, the conditional mean is modelled as

$$g(E[Y|\mathbf{X}]) = \mathbf{X}^\top \boldsymbol{\beta}, \quad \text{for } \boldsymbol{\beta} \in \mathbb{R}^p$$

for some **non-linear link function** $g : \mathcal{M} \rightarrow \mathbb{R}$ (typically continuous, smooth, regular, but nonlinear)

- moreover, the variance structure typically depends on the mean structure

$$\text{Var}[Y|\mathbf{X}] = v(E[Y|\mathbf{X}])$$

where $v : \mathbb{M} \rightarrow (0, \infty)$ is some known (variance) function

Example 1: Logistic regression

□ Logistic regression

- the response variable $Y \in \mathbb{R}$ takes only two possible values, $Y \in \{0, 1\}$
- the conditional distribution of $Y|\mathbf{X}$ is alternative, with $p_x = E[Y|\mathbf{X} = \mathbf{x}]$
- the conditional mean $\mu_x = E[Y|\mathbf{X} = \mathbf{x}]$ is modeled with the **linear predictor** $\mathbf{X}^\top \beta$ using the **logit link function** $g(x) = \log[x/(1-x)]$
- the model assumes the mean structure

$$\text{logit}(\mu_x) = \log \frac{E[Y|\mathbf{X} = \mathbf{x}]}{1 - E[Y|\mathbf{X} = \mathbf{x}]} = \log \frac{P[Y = 1|\mathbf{X} = \mathbf{x}]}{1 - P[Y = 1|\mathbf{X} = \mathbf{x}]} = \mathbf{x}^\top \beta$$

- the model assumes the variance structure which depends on the mean μ_x

$$\text{Var}[Y|\mathbf{X} = \mathbf{x}] = v(\mu_x) = \mu_x(1 - \mu_x)$$

- the model is interpreted in terms of multiplicative comparisons and the parameters are interpreted in terms of the odds ratios

Example 2: Poisson regression

□ Poisson regression

- the response variable $Y \in \mathbb{N} \cup \{0\}$ represents integer counts (including 0)
- the conditional distribution of $Y|\mathbf{X}$ is Poisson, with $\lambda_{\mathbf{x}} = E[Y|\mathbf{X} = \mathbf{x}]$
- the conditional mean $\lambda_{\mathbf{x}} = E[Y|\mathbf{X} = \mathbf{x}]$ is modeled with the **linear predictor** $\mathbf{X}^T \boldsymbol{\beta}$ using the **log link function** $g(x) = \log x$
- the model assumes the mean structure

$$\log(\lambda_{\mathbf{x}}) = \log E[Y|\mathbf{X} = \mathbf{x}] = \mathbf{x}^T \boldsymbol{\beta}$$

- the model assumes the variance structure which depends on the mean $\lambda_{\mathbf{x}} > 0$ and some additional **dispersion** parameter $\phi > 0$

$$\text{Var}[Y|\mathbf{X} = \mathbf{x}] = v(\lambda_{\mathbf{x}})\phi = \phi\lambda_{\mathbf{x}}$$

- the model is interpreted in terms of multiplicative comparisons and the parameters are interpreted in terms of the proportional changes of the conditional expectations

Example 3: Special cases

- ❑ **Classical linear regression model**
 - ❑ continuous response $Y \in \mathbb{R}$
 - ❑ identity link function $g(x) = x$
 - ❑ constant variance function $v(x) = 1$ and $\phi = \sigma^2$

- ❑ **Multinomial regression model**
- ❑ **Exponential data model**
- ❑ ...

2. Nonlinear regression models

- In linear models and generalized linear models as well, the conditional mean is modeled (using a proper link function) as a linear combination of the response variables and the subset of unknown parameters...
- if the class of available models is not reach enough (and we still prefer a parametric model structure) then **nonlinear regression models** can serve a a good alternative...
- the idea in nonlinear models is to use a general parametric (but nonlinear) regression function $f : \mathbb{R}^p \times \mathbb{R}^q \rightarrow \mathbb{R}$, such that

$$E[Y|\mathbf{X}] = f(\mathbf{X}, \beta),$$

where $\mathbf{X} \in \mathbb{R}^p$ and $\beta \in \mathbb{R}^q$

- Note that the nonlinear element (the nonlinear function f) is now introduced on the other side of the classical regression model formula and typically it is not assumed that f should be regular (or continuous, etc.)

Nonlinear regression: Some examples

There are, of course, plenty of different models with various analytical structure and different regularity properties—smoothness, continuity. Typical nonlinear models are, for instance, various population models...

□ Exponential growth model

$$f(x, \beta, \alpha) = \alpha \exp\{X\beta\}$$

→ for some parameters $a > 0$ and $\beta > 0$;

□ Logistic growth model

$$f(X, \beta, \alpha, K) = \frac{K}{1 + be^{-X\beta}};$$

→ for some parameters $\alpha, \beta, K > 0$;

□ Gompertz growth model

$$f(X, \beta, \alpha, K) = K \cdot \exp\{-\beta e^{-\alpha t}\};$$

→ for some parameters $\alpha, \beta, K > 0$;

Solutions for nonlinear regression models

- ❑ Note, that all three nonlinear models above can not be solved by using classical method of the least squares... (no explicit solution can be obtains)
- ❑ Thus, different computation strategies must be used to obtain the model solution—the estimates for the unknown parameters $\alpha, \beta, K > 0$
- ❑ Such computational methods may involve:
 - ❑ reparametrization into a linear model and applying least squares
 - ❑ model approximation and least squares
 - ❑ various iterative solutions
- ❑ Note, that as far as the unknown regression function is unspecified, the corresponding minimization problem may not even be convex!

Generalized nonlinear models

- **Advanced, but still possible....**

$$g(E[Y|\mathbf{X}]) = f(\mathbf{X}, \beta)$$

where two additional sources of nonlinearity are introduced at the same time—the nonlinear link function g and the nonlinear predictor function f

- **Some challenges**
 - mostly, the interpretation of $\beta \in \mathbb{R}$ is not straightforward
 - due to nonlinearity, various computational issues and solution instability
 - difficult statistical inference typically performed by simulations

3. Regression models beyond expectation

□ Least squares

In an ordinary linear regression model (without the normality assumption) the likelihood can not be obtained and the estimates for $\beta' \in \mathbb{R}^p$ are obtained by minimizing least squares

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^p}{\text{Argmin}} \sum_{i=1}^N (Y_i - \mathbf{x}_i^\top \beta)^2$$

3. Regression models beyond expectation

□ Least squares

In an ordinary linear regression model (without the normality assumption) the likelihood can not be obtained and the estimates for $\beta' \in \mathbb{R}^p$ are obtained by minimizing least squares

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^p}{\text{Argmin}} \sum_{i=1}^N (Y_i - \mathbf{x}_i^\top \beta)^2$$

□ Maximum likelihood

In a normal linear regression model (under the normality assumption) the full likelihood for $\beta \in \mathbb{R}^p$ and $\sigma > 0$ can be formulated and the estimates are obtained by maximizing the likelihood function

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^p; \sigma^2 > 0}{\text{Argmax}} (2\pi\sigma^2)^{-N/2} \cdot \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N (Y_i - \mathbf{x}_i^\top \beta)^2 \right\}$$

Expectation and beyond

- For some real random variable $X \sim F_X$ (and the density f with respect to the Lebesgue or count measure) and some measurable function $h : \mathbb{R} \rightarrow \mathbb{R}$ we can obtain the expectation (if the integral exists) as

$$Eh(X) = \int_{\mathbb{R}} h(x) dF_X(x) = \int_{\mathbb{R}} h(x) f(x) dx$$

Expectation and beyond

- For some real random variable $X \sim F_X$ (and the density f with respect to the Lebesgue or count measure) and some measurable function $h : \mathbb{R} \rightarrow \mathbb{R}$ we can obtain the expectation (if the integral exists) as

$$Eh(X) = \int_{\mathbb{R}} h(x)dF_X(x) = \int_{\mathbb{R}} h(x)f(x)dx$$

- For the random sample X_1, \dots, X_N drawn from the same distribution as the distribution of $X \sim F_X$ we can construct the empirical distribution function F_N and the empirical counterpart for $Eh(X)$ (i.e., the empirical estimate)

$$\widehat{Eh(X)} = \int_{\mathbb{R}} h(x)dF_N(x) = \sum_{i=1}^N h(X_i)$$

Expectation and beyond

- For some real random variable $X \sim F_X$ (and the density f with respect to the Lebesgue or count measure) and some measurable function $h : \mathbb{R} \rightarrow \mathbb{R}$ we can obtain the expectation (if the integral exists) as

$$Eh(X) = \int_{\mathbb{R}} h(x) dF_X(x) = \int_{\mathbb{R}} h(x) f(x) dx$$

- For the random sample X_1, \dots, X_N drawn from the same distribution as the distribution of $X \sim F_X$ we can construct the empirical distribution function F_N and the empirical counterpart for $Eh(X)$ (i.e., the empirical estimate)

$$\widehat{Eh(X)} = \int_{\mathbb{R}} h(x) dF_N(x) = \sum_{i=1}^N h(X_i)$$

- The quantity (parameter) $\mu_h = Eh(X)$ is sometimes called the theoretical is called the **theoretical functional of the distribution F_X** while the quantity $\widehat{\mu}_h = \widehat{Eh(X)}$ is called the **(empirical) functional of the empirical distribution F_N** (different functions can be used in place of h)

Some common choices of the function h

- the expectation (theoretical quantity EX) and the average (empirical quantity \bar{X}_N) can be both obtained by a minimization problem with the choice of $h(x) = (x - a)^2$ where

- $EX = \mathop{\text{Argmin}}_{a \in \mathbb{R}} E(X - a)^2 = \mathop{\text{Argmin}}_{a \in \mathbb{R}} \int_{\mathbb{R}} (x - a)^2 dF_X(x)$

- $\bar{X}_N = \widehat{EX} = \mathop{\text{Argmin}}_{a \in \mathbb{R}} \int_{\mathbb{R}} (x - a)^2 dF_N(x) = \mathop{\text{Argmin}}_{a \in \mathbb{R}} \sum_{i=1}^N (X_i - a)^2$

Some common choices of the function h

- the expectation (theoretical quantity EX) and the average (empirical quantity \bar{X}_N) can be both obtained by a minimization problem with the choice of $h(x) = (x - a)^2$ where

- $EX = \mathop{\text{Argmin}}_{a \in \mathbb{R}} E(X - a)^2 = \mathop{\text{Argmin}}_{a \in \mathbb{R}} \int_{\mathbb{R}} (x - a)^2 dF_X(x)$

- $\bar{X}_N = \widehat{EX} = \mathop{\text{Argmin}}_{a \in \mathbb{R}} \int_{\mathbb{R}} (x - a)^2 dF_N(x) = \mathop{\text{Argmin}}_{a \in \mathbb{R}} \sum_{i=1}^N (X_i - a)^2$

- Note that in both cases we actually formulate the **least squares problem** (theoretical and empirical) and the solution is the theoretical mean and the empirical average (i.e., the estimate for the mean)

This principle can be generalized even further—for the regression concepts and different forms of the function h

Some common choices of the function h

- the expectation (theoretical quantity EX) and the average (empirical quantity \bar{X}_N) can be both obtained by a minimization problem with the choice of $h(x) = (x - a)^2$ where

$$\square EX = \mathop{\text{Argmin}}_{a \in \mathbb{R}} E(X - a)^2 = \mathop{\text{Argmin}}_{a \in \mathbb{R}} \int_{\mathbb{R}} (x - a)^2 dF_X(x)$$

$$\square \bar{X}_N = \widehat{EX} = \mathop{\text{Argmin}}_{a \in \mathbb{R}} \int_{\mathbb{R}} (x - a)^2 dF_N(x) = \mathop{\text{Argmin}}_{a \in \mathbb{R}} \sum_{i=1}^N (X_i - a)^2$$

- Note that in both cases we actually formulate the **least squares problem** (theoretical and empirical) and the solution is the theoretical mean and the empirical average (i.e., the estimate for the mean)

This principle can be generalized even further—for the regression concepts and different forms of the function h

- typical choices for h include: median regression for $h(x) = |x|$; quantile regression for $h_\tau(x) = \tau(x - \mathbb{I}_{\{x < 0\}})$; expectile regression for $h_\tau(x) = |\tau - \mathbb{I}_{\{x < 0\}}|x^2$; robust regression for $h(x) = \rho(x)$

Basic properties of the regression variants

□ Median regression

- more robust than the standard least squares regression
- for symmetric error distributions the median corresponds with the mean
- easy and straightforward interpretation of the estimated parameters

□ Quantile regression

- generalization of the median regression (which is obtained for $\tau = 0.5$)
- provides a complex insight about the conditional distribution of $Y|X$
- relatively easy interpretation but not that much popular in practice

□ Expectile regression

- generalization of the least squares (which are obtained for $\tau = 0.5$)
- expectiles form elastic and coherent risk measures (unlike quantiles)
- relatively difficult interpretation but very popular in risk theory

□ Robust regression

- generalization of the regression for outliers and heavy-tailed distributions
- least squares for $\rho(x) = X^2$; median regression for $\rho(x) = |x|$; maximum likelihood for $\rho(x) = -\log(x)$
- other choices are common in practice as well (e.g., Huber function, Tukey function, Andrew's function, ...)

Exam terms

- ❑ **Utorok, 20.05.2025** (starting at 10:40 in K4)
- ❑ **Štvrtok, 22.05.2025** (starting at 12:20 in Praktikum KPMS)
- ❑ **Utorok, 27.05.2025** (starting at 9:00 in K5) (!!)
- ❑ **Štvrtok, 29.05.2025** (starting at 9:00 in K11)
- ❑ **Ponedelok, 02.06.2025** (starting at 9:00 in K11)
- ❑ **Ponedelok, 09.06.2025** (starting at 9:00 in K11)
- ❑ **Štvrtok, 26.06.2025** (starting at 9:00 in K11)

- ❑ At least one other exam term in **September 2025**