# Why linear models?

Roger Levy

Massachusetts Institute of Technology

March 10, 2025

# Why linear models?

For modeling conditional distributions $P(Y|X)$, we get to linear regression from the following assumptions:

1. The expected response $E[Y|X]$ is a <span style="color:red">linear combination</span> of the predictors:

$$\widehat{y} = \beta \cdot \mathsf{x}$$

# Why linear models?

For modeling conditional distributions $P(Y|X)$, we get to linear regression from the following assumptions:

1. The expected response $E[Y|X]$ is a linear combination of the predictors:

$$\widehat{y} = \beta \cdot \mathsf{x}$$

2. Residual error, or noise, is distributed normally with mean zero around the expected response:

$$y \sim \mathcal{N}(\widehat{y}, \sigma^2)$$

# Why linear models?

For modeling conditional distributions $P(Y|X)$, we get to linear regression from the following assumptions:

1. The expected response $E[Y|X]$ is a linear combination of the predictors:

$$\widehat{y} = \beta \cdot \mathsf{x}$$

2. Residual error, or noise, is distributed normally with mean zero around the expected response:

$$y \sim \mathcal{N}(\widehat{y}, \sigma^2)$$

# Why linear models?

For modeling conditional distributions $P(Y|X)$, we get to linear regression from the following assumptions:

1. The expected response $E[Y|X]$ is a linear combination of the predictors:

$$\widehat{y} = \beta \cdot \mathsf{x}$$

2. Residual error, or noise, is distributed normally with mean zero around the expected response:

$$y \sim \mathcal{N}(\widehat{y}, \sigma^2)$$

But, why linear regression?

# Linear effect of predictors

$$\boxed{\widehat{y} = \beta \cdot \mathsf{x}}$$
XXX

$$y \sim \mathcal{N}(\widehat{y}, \sigma^2)$$

Three major reasons why this is a good assumption to use:

# Normally distributed residual error

$$y \sim \mathcal{N}(\widehat{y}, \sigma^2)$$

Three major reasons why this is a good assumption to use:

- ▶ It has good computational properties

# Normally distributed residual error

$$y \sim \mathcal{N}(\widehat{y}, \sigma^2)$$

Three major reasons why this is a good assumption to use:

- It has good computational properties
- There are *a priori* mathematical reasons to expect residual error to often look Gaussian(-ish)

# Normally distributed residual error

$$y \sim \mathcal{N}(\widehat{y}, \sigma^2)$$

Three major reasons why this is a good assumption to use:

- It has good computational properties
- There are *a priori* mathematical reasons to expect residual error to often look Gaussian(-ish)
- It is robust to mild-to-moderate violations

# Good computational properties of Gaussian-error assumption

$$P(y|\widehat{y}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y - \widehat{y})^2}{2\sigma^2}\right]$$

▶ Regardless of the value of $\sigma^2$, this quantity is minimized when $\sum_i (y_i - \widehat{y}_i)^2$ is minimized across the dataset

# Good computational properties of Gaussian-error assumption

$$P(y|\hat{y}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y - \hat{y})^2}{2\sigma^2}\right]$$

▶ Regardless of the value of $\sigma^2$, this quantity is minimized when $\sum_i (y_i - \hat{y}_i)^2$ is minimized across the dataset

▶ Therefore, maximizing data likelihood is equivalent to **minimizing the sum of squared residual errors**

# Good computational properties of Gaussian-error assumption

$$P(y|\widehat{y}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y - \widehat{y})^2}{2\sigma^2}\right]$$

▶ Regardless of the value of $\sigma^2$, this quantity is minimized when $\sum_i (y_i - \widehat{y}_i)^2$ is minimized across the dataset

▶ Therefore, maximizing data likelihood is equivalent to **minimizing the sum of squared residual errors**

    ▶ This gives rise to the name commonly used for linear regression, ordinary least squares

# The least-squares fit for linear regression

Let $\text{RSS}(\beta) = \sum_i (y_i - \widehat{y}_i)^2$ and use matrix notation $\mathbf{y} = \mathsf{X}\beta + \epsilon$:

$$\text{RSS}(\beta) = (\mathbf{y} - \mathsf{X}\beta)^\top (\mathbf{y} - \mathsf{X}\beta)$$

This is minimized when its gradient with respect to $\beta$ is 0, i.e., $\nabla_\beta(\text{RSS}(\beta)) = 0$. Using vector calculus we derive:

$$
\begin{aligned}
0 &= \nabla_\beta \left( (\mathbf{y} - \mathsf{X}\beta)^\top (\mathbf{y} - \mathsf{X}\beta) \right) \\
&= \left( \nabla_\beta (y - X\beta) \right)^\top (y - X\beta) + \left( \nabla_\beta (y - X\beta) \right)(y - X\beta)^\top \\
&= -\mathsf{X}^\top (y - X\beta) - \mathsf{X}(y - X\beta)^\top \\
&= -2 \left( \mathsf{X}^\top \mathbf{y} - \mathsf{X}^\top \mathsf{X}\beta \right)
\end{aligned}
$$

If $\mathsf{X}^\top \mathsf{X}$ is invertible, then this equation is solved when:

$$\beta = \left( \mathsf{X}^\top \mathsf{X} \right)^{-1} \mathsf{X}^\top \mathbf{y}$$

This is thus the OLS, and equivalently the ML, estimate $\widehat{\beta}$.

▶ In regression modeling with a continuous response, a frequently plausible working assumption is that the residual error involves the **additive influence** of **large numbers** of **separate** (unconditionally independent) **stochastic factors** that are **not captured by the predictors included in the regression model**.

▶ In regression modeling with a continuous response, a frequently plausible working assumption is that the residual error involves the **additive influence** of **large numbers** of **separate** (unconditionally independent) **stochastic factors** that are **not captured by the predictors included in the regression model**.

▶ In this setting, we can invoke the Central Limit Theorem, stated here somewhat informally:

> The mean of $n$ independently distributed random variables approaches a normal distribution as $n \to \infty$.

# How bad is it if residual error is not truly Gaussian?

Even if residual error is not truly Gaussian, OLS linear regression still has desirable properties:

▶ So long as errors are IID, OLS is the Best Linear Unbiased Estimator (BLUE) of the regression coefficients—this is the Gauss–Markov theorem

Even if residual error is not truly Gaussian, OLS linear regression still has desirable properties:

- ▶ So long as errors are IID, OLS is the Best Linear Unbiased Estimator (BLUE) of the regression coefficients—this is the Gauss–Markov theorem
  - ▶ ("Best" here means lowest-variance)

# How bad is it if residual error is not truly Gaussian?

Even if residual error is not truly Gaussian, OLS linear regression still has desirable properties:

- ▶ So long as errors are IID, OLS is the Best Linear Unbiased Estimator (BLUE) of the regression coefficients—this is the Gauss–Markov theorem
    - ▶ ("Best" here means lowest-variance)

# How bad is it if residual error is not truly Gaussian?

Even if residual error is not truly Gaussian, OLS linear regression still has desirable properties:

- So long as errors are IID, OLS is the Best Linear Unbiased Estimator (BLUE) of the regression coefficients—this is the Gauss–Markov theorem
    - ("Best" here means lowest-variance)

However, non-normal residual error affects correctness of inferences regarding coefficient estimates, beyond simple unbiasedness. We now turn to a large family of techniques that manages the resulting issues.

# References I