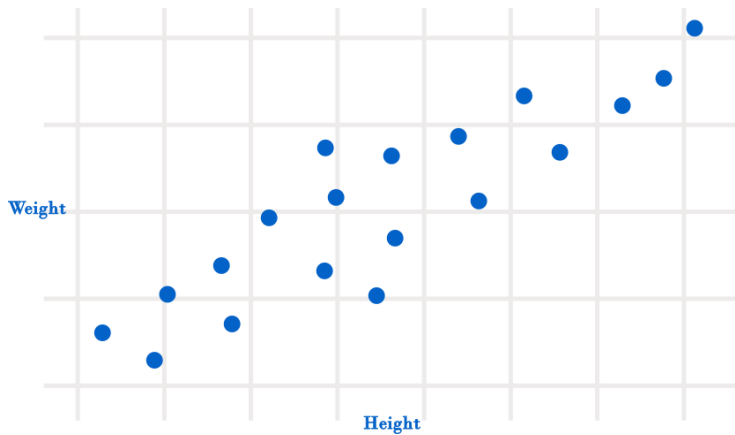


# Linear Regression

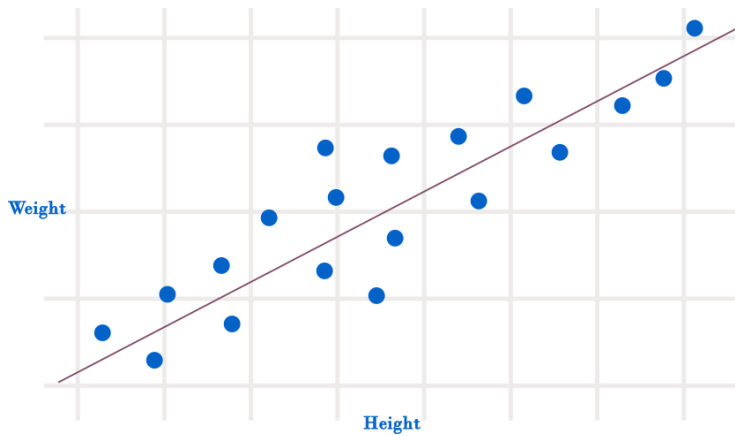
By Van Dinh Tran

January 27, 2025

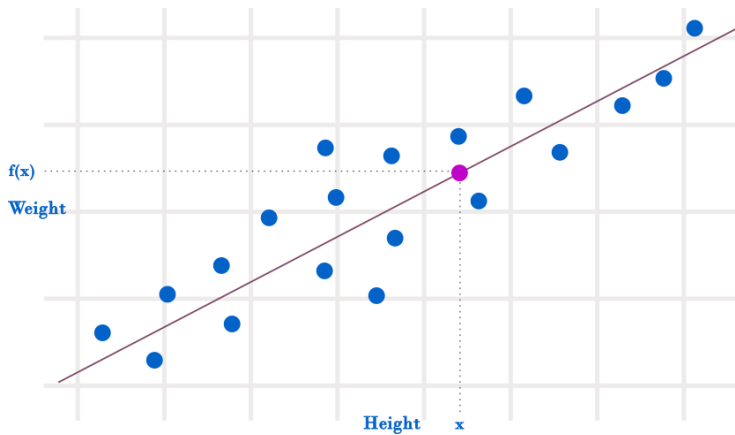
# Motivation



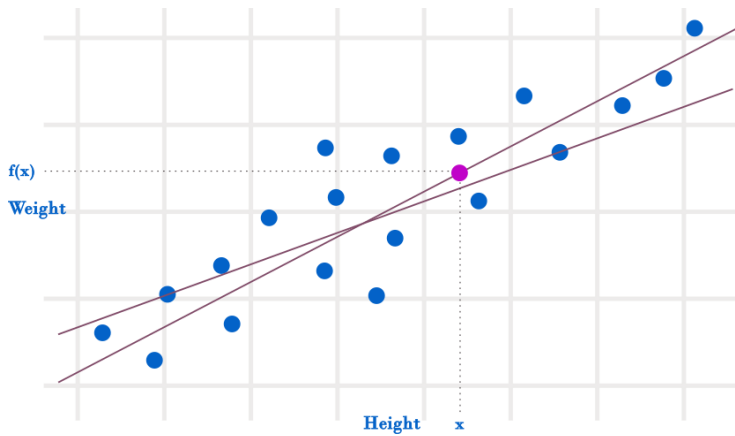
# Motivation



# Motivation



# Motivation



# Linear regression

- **Data:**

$$D = \{(x_i, y_i)\} \mid x_i \in X \subseteq \mathbb{R}^d, y_i \in Y \subseteq \mathbb{R}\}, |D| = N$$

$x_i$ : input vector or instance;  $y_i$ : target (label)

- **Hypothesis space:**

$$\mathcal{H} = \{h_w(x) \mid h_w(x) = \mathbf{w}^\top x; h_w(x) : X \longrightarrow Y\}$$

( $\mathcal{H}$  is a set of linear functions.)

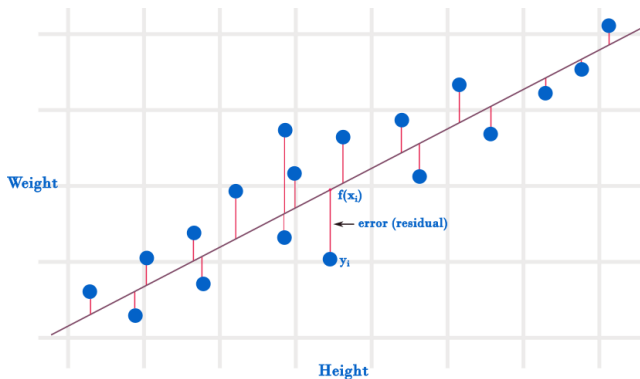
- **Loss function:**

$$L(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \text{loss}(h_w(x_i), y_i)$$

- **Optimization:** Find  $\mathbf{w}$  (hypothesis) that minimizes  $L(\mathbf{w})$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} L(\mathbf{w})$$

# Ordinary Least Square (OLS) regression



Loss function:

$$L(w) = \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2$$

$$w^* = \arg \min_w \left[ \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2 \right]$$

# Closed form solution

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left[ \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 \right]$$

**Idea:** Find  $\mathbf{w}^*$  by setting partial derivatives to zero.

$$\frac{\partial L}{\partial w_j} = \frac{2}{N} \sum_{i=1}^N x_{ij} [\mathbf{w}^\top \mathbf{x}_i - y_i] = \frac{2}{N} \sum_{i=1}^N x_{ij} \left[ \sum_{k=1}^d w_k x_{ik} - y_i \right]$$

- Let  $\frac{\partial L}{\partial w_j} = 0$ :

$$\sum_{k=1}^d \left[ \sum_{i=1}^N x_{ij} x_{ik} \right] w_k - \sum_{i=1}^N x_{ij} y_i = 0$$

- Notice that It has the form of  $\mathbf{A}\mathbf{w} - \mathbf{b} = 0$  with

$$A_{ij} = \sum_{i=1}^N x_{ij} x_{ik}; \quad b_i = \sum_{i=1}^N x_{ij} y_i$$

- Solving this linear system, we get  $\mathbf{w}^*$

Note that this solution is particular for this linear regression, not general linear regression.



## Closed form solution: Matrix form

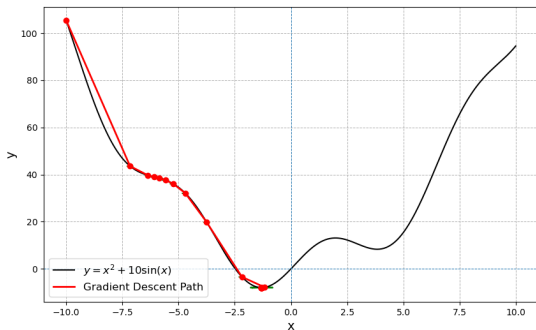
$$\begin{aligned}L(\mathbf{w}) &= \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \\&= (\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) \\&= (\mathbf{w}^\top \mathbf{X}^\top - \mathbf{y}^\top)(\mathbf{X}\mathbf{w} - \mathbf{y}) \\&= \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} \mathbf{w} + \mathbf{y}^\top \mathbf{y}\end{aligned}$$

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{w}} &= 2\mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{y} \\&= 2\mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{X}^\top \mathbf{y}\end{aligned}$$

$$2\mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{X}^\top \mathbf{y} = 0 \implies \mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

- Utilizing SVD
- Note that  $\mathbf{X}^\top \mathbf{X}$  is invertible if and only if it has full rank
- We can use gradient descent to solve

# Gradient Descent



- Gradient descent is an iterative optimization algorithm to find a function's optimum.
- Commonly used in machine learning to optimize the parameters of complex models.
- Based on the property of gradient that it shows the direction of the steepest ascent of a function

# Gradient Descent

$$\min_w L(w)$$

## Gradient Descent:

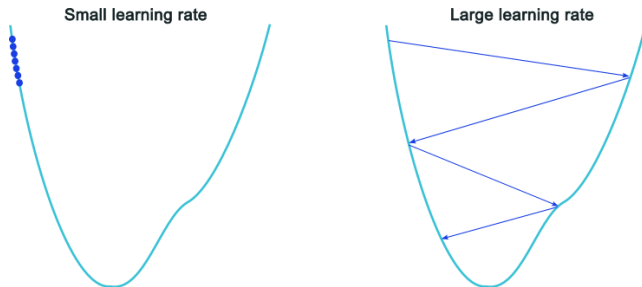
- Pick some value for  $w$
- Iteratively update  $w$  until convergence:

$$w \leftarrow w - \alpha * \nabla L$$

where  $\alpha$  is the learning rate and  $\nabla L$  is defined as follows:

$$\nabla L = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \\ \vdots \\ \frac{\partial L}{\partial w_d} \end{bmatrix}$$

# Gradient Descent: learning rate



- Learning rate ( $\alpha$ ): a non-negative hyperparameter
- $\alpha$  is small: slow convergence, high precise convergence, and risk of getting stuck at local optima
- $\alpha$  is large: fast convergence, overshooting optimum, and instability

# OLS: pros and cons

- **Pros:**

- ▶ Simplicity
- ▶ Interpretability
- ▶ Efficiency

- **Cons:**

- ▶ Linearity assumption
- ▶ Sensitivity to outliers
- ▶ Multicollinearity
- ▶ No Automatic Feature Selection
- ▶ Prone to overfitting
- ▶ Assumption of Independence
- ▶ Limited to Linear Relationships

# OLS feature importance

- Standardization: Standardizing the predictors allows for direct comparison of coefficients
- t-statistics, p-value
- Adjusted R-squared contribution
- Combination of different methods

Note: “*statsmodels*” is a good Python package used in practice to interpret the importance of features.

# Regression evaluation metrics

- Mean absolute error (MAE):

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$

- Mean square error (MSE):

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

- Root mean square error (RMSE):

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}$$

- $R^2$  error:

$$R^2 = 1 - \frac{\text{SS}_{\text{res}}}{\text{SS}_{\text{tot}}}$$

- Adjusted  $R^2$  error:

$$R^2 = 1 - \frac{(1-R^2)(n-1)}{N-d-1}$$

# Regularized linear regression

- Regularization refers to techniques used to calibrate machine learning models to minimize the adjusted loss function and prevent over-fitting or under-fitting.
- It is most commonly known for its use against over-fitting
- Linear regression with regularization:
  - Lasso: using L1 norm ( $\|\cdot\|_1$ )
  - Ridge: using L2 norm ( $\|\cdot\|_2$ )
  - Elastic Net: combining L1 and L2 ( $\|\cdot\|_1$  and  $\alpha(\|\cdot\|_2)$ )
- Norms of  $w = (w_1, w_2, \dots, w_d) \in \mathbb{R}^d$ :
  - Norm 1:  $\|w\|_1 = |w_1| + |w_2| + \dots + |w_d|$
  - Norm 2:  $\|w\|_2 = \sqrt{w_1^2 + w_2^2 + \dots + w_d^2}$
  - Norm infinity:  $\|w\|_\infty = \text{Max}(|w_1|, |w_2|, \dots, |w_d|)$



# Ridge regression

- Ridge aims to shrink linear regression coefficients, i.e. parameters are pushed toward zero.
- It adds constraints to the objective function of the simple linear regression using the L2 norm.

$$L(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

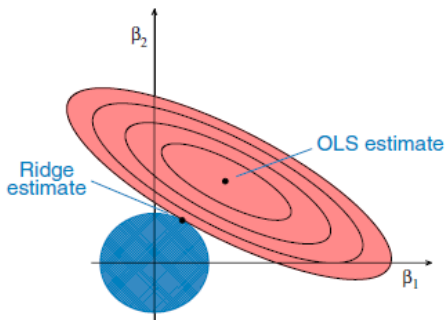
( $\|\mathbf{w}\|_2^2$ : regularization term,  $\lambda$ : regularization strength hyperparameter)

$$\frac{\partial L}{\partial \mathbf{w}} = 2\mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{X}^\top \mathbf{y} + 2\lambda \mathbf{w} = 0$$

$$\implies \mathbf{w}^* = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

Note that  $(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})$  is positive-definite, so is invertable.

# Ridge regression



Geometric interpretation of Ridge regression

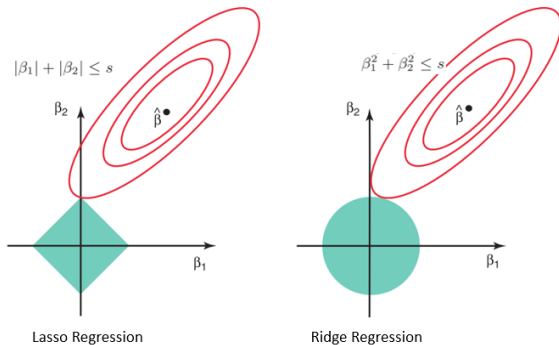
# Lasso regression

- Like Ridge regression, it is a regularization technique that uses shrinkage
- The lasso procedure encourages simple, sparse models, i.e. models with fewer parameters. Thus, it has a feature selection capability
- It adds constraints to the objective function of the simple linear regression using the L1 norm.

$$L(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1$$

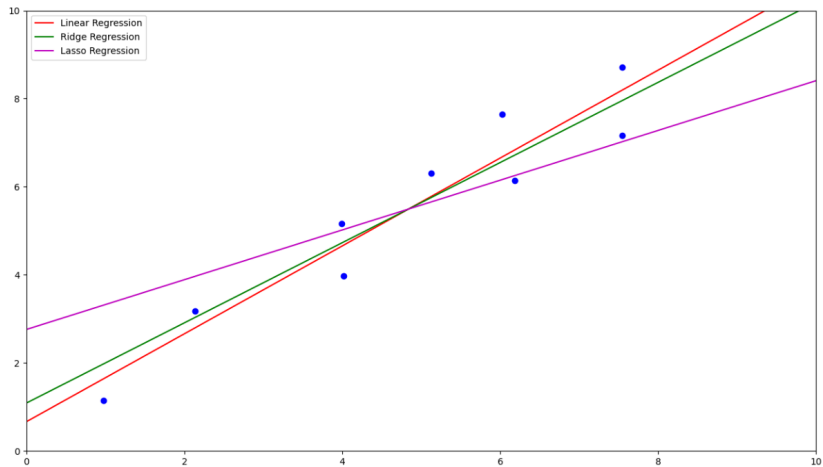
- $L(\mathbf{w})$  is not differentiable at zero ( $\mathbf{w}=0$ ).
- Subgradient can be used to calculate zero gradients.
- Subgradient descent is used instead of gradient descent.

# Lasso regression



Ridge regression shrinks coefficients towards zero, while Lasso tends to give a subset of coefficients zero, leading to a sparse solution.

# OLS, Ridge and Lasso



OLS, Ridge ( $\lambda = 4$ ), Lasso ( $\lambda = 2$ )

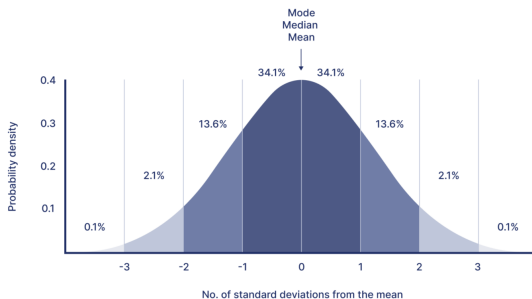
# Review Normal distribution (Gaussian distribution)

- Univariate:

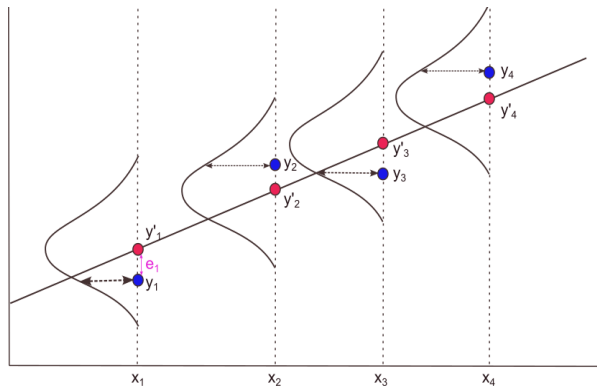
- $\mathbf{x} \sim \mathcal{N}(\mu, \sigma^2)$
- Pdf:  $f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$

- Multivariate:

- $\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)$
- Pdf:  $f(\mathbf{x}) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^\top \Sigma^{-1}(\mathbf{x} - \mu)\right)$



# Maximum likelihood regression



- Blue points are true values (labels)
- Red points are predicted values
- For each  $x_i$ , we estimate a Gaussian distribution over  $y_i$  values,  $N(w^T x_i, \sigma)$

# Maximum likelihood regression

- Assuming that  $y_i \sim \mathcal{N}(w^\top x_i, \sigma^2)$  and  $y_i$  are independent
- The likelihood of observing  $\{y_1, y_2, \dots, y_N\}$  is as follow:

$$\mathcal{L}(w, \sigma) = p(y|X, w, \sigma) = \prod_{i=1}^N p(y_i|x_i, w, \sigma)$$

$$\begin{aligned}\mathcal{L}(w, \sigma) &= \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - w^\top x_i)^2}{2\sigma^2}} \\ &= \frac{1}{(\sqrt{2\pi\sigma^2})^N} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - w^\top x_i)^2} \\ &= \frac{1}{(\sqrt{2\pi\sigma^2})^N} e^{-\frac{1}{2\sigma^2} (y - Xw)^\top (y - Xw)} \quad (\text{Vector form})\end{aligned}$$

- Note that maximizing  $\mathcal{L}(w, \sigma)$  is equivalent to minimizing the loss function,  $L$ , of linear regression.

$$\max_{w, \sigma} \mathcal{L}(w, \sigma) \iff \min_w L(w)$$



# Maximum likelihood regression

$$\max_{\mathbf{w}, \sigma} \mathcal{L}(\mathbf{w}, \sigma) \iff \max_{\mathbf{w}, \sigma} \text{Log}(\mathcal{L}(\mathbf{w}, \sigma))$$

$$\begin{aligned}\text{Log}(\mathcal{L}(\mathbf{w}, \sigma)) &= \text{Log}\left(\frac{1}{(\sqrt{2\pi\sigma^2})^N} e^{-\frac{1}{2\sigma^2} (\mathbf{y}-\mathbf{X}\mathbf{w})^\top (\mathbf{y}-\mathbf{X}\mathbf{w})}\right) \\&= \text{Log}\left(\frac{1}{(\sqrt{2\pi\sigma^2})^N}\right) + \text{Log}(e^{-\frac{1}{2\sigma^2} (\mathbf{y}-\mathbf{X}\mathbf{w})^\top (\mathbf{y}-\mathbf{X}\mathbf{w})}) \\&= -\frac{N}{2} \text{Log}(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) \\&= -\frac{N}{2} \text{Log}(2\pi\sigma^2) - \frac{1}{2\sigma^2} [(\mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}\mathbf{w} - \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X}\mathbf{w})] \\&= -\frac{N}{2} \text{Log}(2\pi\sigma^2) - \frac{1}{2\sigma^2} [(\mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X}\mathbf{w} + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X}\mathbf{w})]\end{aligned}$$

# Maximum likelihood regression

Note that  $\frac{\partial A\mathbf{w}}{\partial \mathbf{w}} = A^\top$  and  $\frac{\partial \mathbf{w}^\top A\mathbf{w}}{\partial \mathbf{w}} = 2A^\top \mathbf{w}$  (Matrix differentiation)

$$\frac{\partial \text{Log}(\mathcal{L}(\mathbf{w}, \sigma))}{\partial \mathbf{w}} = 0 - 2X^\top \mathbf{y} + 2XX^\top \mathbf{w}$$

$$\frac{\partial \text{Log}(\mathcal{L}(\mathbf{w}, \sigma))}{\partial \mathbf{w}} = 0$$

$$\implies \mathbf{w} = (XX^\top)^{-1}X^\top \mathbf{y}$$

$$\frac{\partial \text{Log}(\mathcal{L}(\mathbf{w}, \sigma))}{\partial \sigma} = -\frac{N}{\sigma} + \frac{1}{\sigma^3}(\mathbf{y} - X\mathbf{w})^2$$

$$\frac{\partial \text{Log}(\mathcal{L}(\mathbf{w}, \sigma))}{\partial \sigma} = 0$$

$$\implies \sigma^2 = \frac{1}{N}(\mathbf{y} - X\mathbf{w})^2$$

**Prediction:**  $\hat{\mathbf{y}} = \hat{X}\mathbf{w}$

# Maximum likelihood regression

The key advantage of maximum likelihood regression:

- Variance of estimators
- Handling different distributions: maximum likelihood regression can handle different error distributions, while other linear regression models assume normally distributed errors.
- Complex models: It can easily be extended to handle more complex models, including those with non-linear relationships