# Outline

- Input representation

- Linear Classification

- Linear Regression     **regression $\equiv$ real-valued output**

- Nonlinear Transformation

# Credit again

**Classification**: Credit approval  (yes/no)

**Regression**: Credit line  (dollar amount)

Input:  $\mathbf{x} =$

| age | 23 years |
|---|---|
| annual salary | $30,000 |
| years in residence | 1 year |
| years in job | 1 year |
| current debt | $15,000 |
| ... | ... |

Linear regression output:  $h(\mathbf{x}) = \sum\limits_{i=0}^{d} w_i \, x_i = \mathbf{w}^\mathsf{T}\mathbf{x}$

# The data set

Credit officers decide on credit lines:

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \cdots , (\mathbf{x}_N, y_N)$$

$y_n \in \mathbb{R}$ is the credit line for customer $\mathbf{x}_n$.
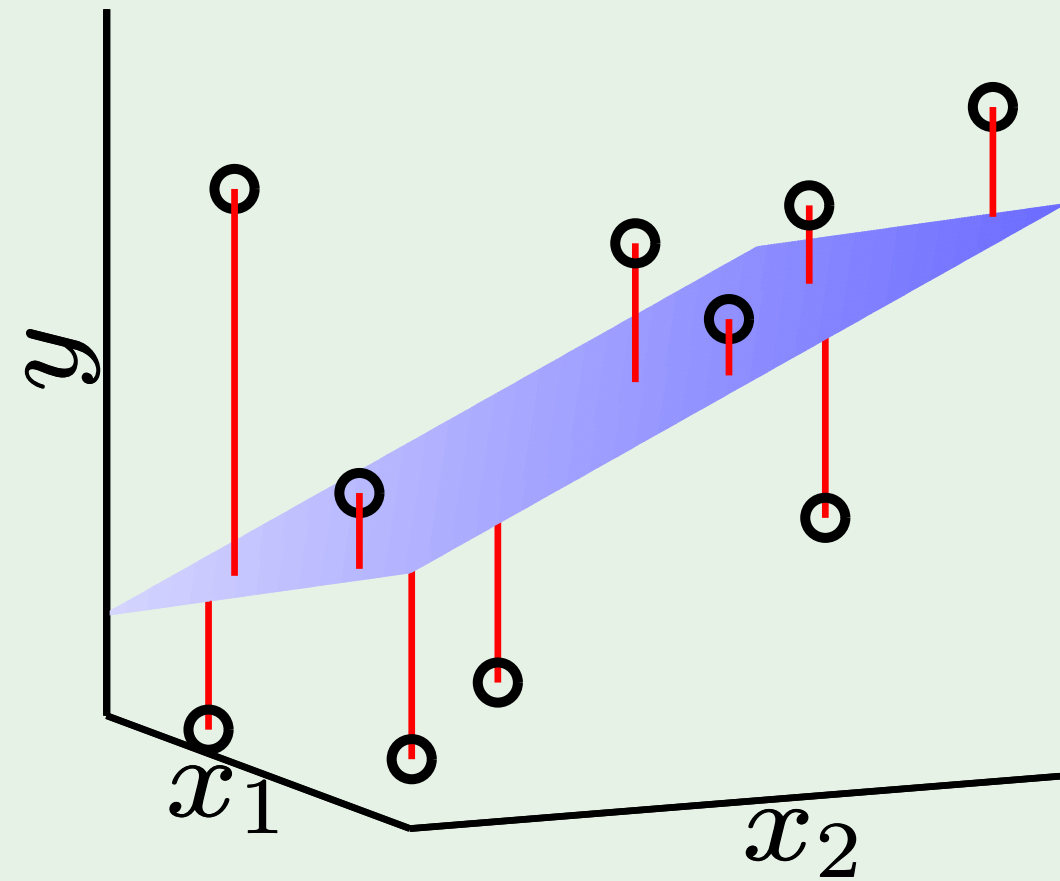
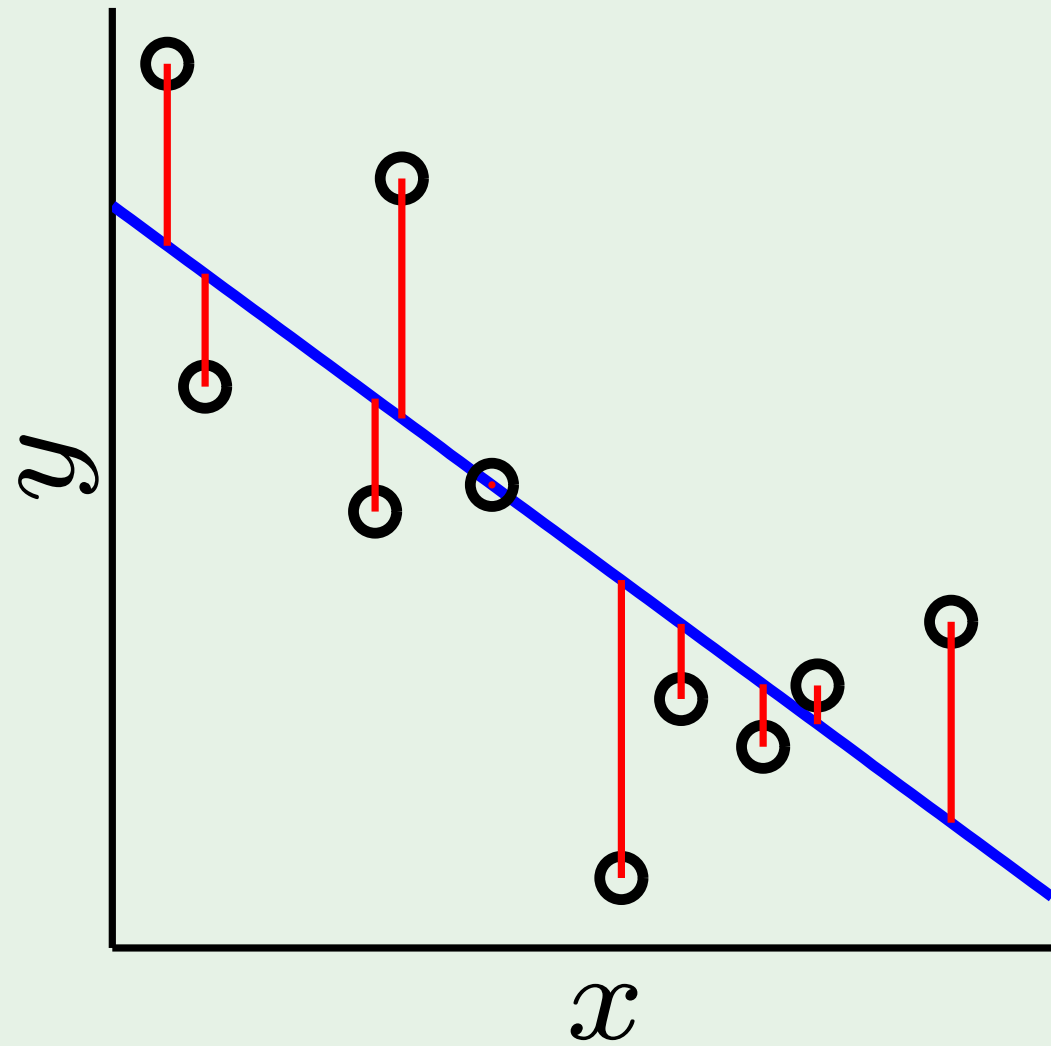Linear regression tries to replicate that.

# How to measure the error

How well does $h(\mathbf{x}) = \mathbf{w}^{\mathsf{T}}\mathbf{x}$ approximate $f(\mathbf{x})$?

In linear regression, we use squared error $(h(\mathbf{x}) - f(\mathbf{x}))^2$

$$\text{in-sample error: } E_{\text{in}}(h) = \frac{1}{N}\sum_{n=1}^{N}(h(\mathbf{x}_n) - y_n)^2$$

# Illustration of linear regression

# The expression for $E_{\mathsf{in}}$

$$E_{\mathsf{in}}(\mathbf{w}) \;=\; \frac{1}{N}\sum_{n=1}^{N}\left(\mathbf{w}^{\mathsf{T}}\mathbf{x}_n - y_n\right)^2$$

$$=\; \frac{1}{N}\|\mathrm{X}\mathbf{w} - \mathbf{y}\|^2$$

where $\quad \mathrm{X} = \begin{bmatrix} -\mathbf{x}_1^{\mathsf{T}}- \\ -\mathbf{x}_2^{\mathsf{T}}- \\ \vdots \\ -\mathbf{x}_N^{\mathsf{T}}- \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$

# Minimizing $E_{\mathsf{in}}$

$$E_{\mathsf{in}}(\mathbf{w}) = \tfrac{1}{N}\|\mathrm{X}\mathbf{w} - \mathbf{y}\|^2$$

$$\nabla E_{\mathsf{in}}(\mathbf{w}) = \tfrac{2}{N}\mathrm{X}^{\top}(\mathrm{X}\mathbf{w} - \mathbf{y}) = \mathbf{0}$$

$$\mathrm{X}^{\top}\mathrm{X}\mathbf{w} = \mathrm{X}^{\top}\mathbf{y}$$

$$\mathbf{w} = \mathrm{X}^{\dagger}\mathbf{y} \;\; \text{where} \;\; \mathrm{X}^{\dagger} = (\mathrm{X}^{\top}\mathrm{X})^{-1}\mathrm{X}^{\top}$$

$$\mathrm{X}^{\dagger} \text{ is the 'pseudo-inverse' of } \mathrm{X}$$

# The pseudo-inverse

$$\mathrm{X}^{\dagger} = (\mathrm{X}^{\mathsf{T}}\mathrm{X})^{-1}\mathrm{X}^{\mathsf{T}}$$

$$\left( \underbrace{\begin{bmatrix} \phantom{xx} \end{bmatrix}}_{d+1 \ \times \ d+1} \right)^{-1} \underbrace{\begin{bmatrix} \phantom{xxxxxx} \end{bmatrix}}_{d+1 \ \times \ N}$$

$$\underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxx}}_{d+1 \ \times \ N}$$

# The linear regression algorithm

1: Construct the matrix $X$ and the vector $\mathbf{y}$ from the data set $(\mathbf{x}_1, y_1), \cdots, (\mathbf{x}_N, y_N)$ as follows

$$X = \underbrace{\begin{bmatrix} -\mathbf{x}_1^\top- \\ -\mathbf{x}_2^\top- \\ \vdots \\ -\mathbf{x}_N^\top- \end{bmatrix}}_{\text{input data matrix}}, \qquad \mathbf{y} = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\text{target vector}}.$$

2: Compute the pseudo-inverse $X^\dagger = (X^\top X)^{-1} X^\top$.

3: Return $\mathbf{w} = X^\dagger \mathbf{y}$.

# Linear regression for classification

Linear regression learns a real-valued function $y = f(\mathbf{x}) \in \mathbb{R}$

Binary-valued functions are also real-valued! $\pm 1 \in \mathbb{R}$

Use linear regression to get $\mathbf{w}$ where $\color{red}{\mathbf{w}^{\mathsf{T}}\mathbf{x}_n} \approx \color{blue}{y_n} = \pm 1$

In this case, $\color{red}{\text{sign}(\mathbf{w}^{\mathsf{T}}\mathbf{x}_n)}$ is likely to agree with $\color{blue}{y_n} = \pm 1$

Good initial weights for classification

# Linear regression boundary