# Least Squres Objective $\hat{y_i} = wx_i$ linear regression

- Instead of "exact $y_i$", we evaluate the size of error in prediction
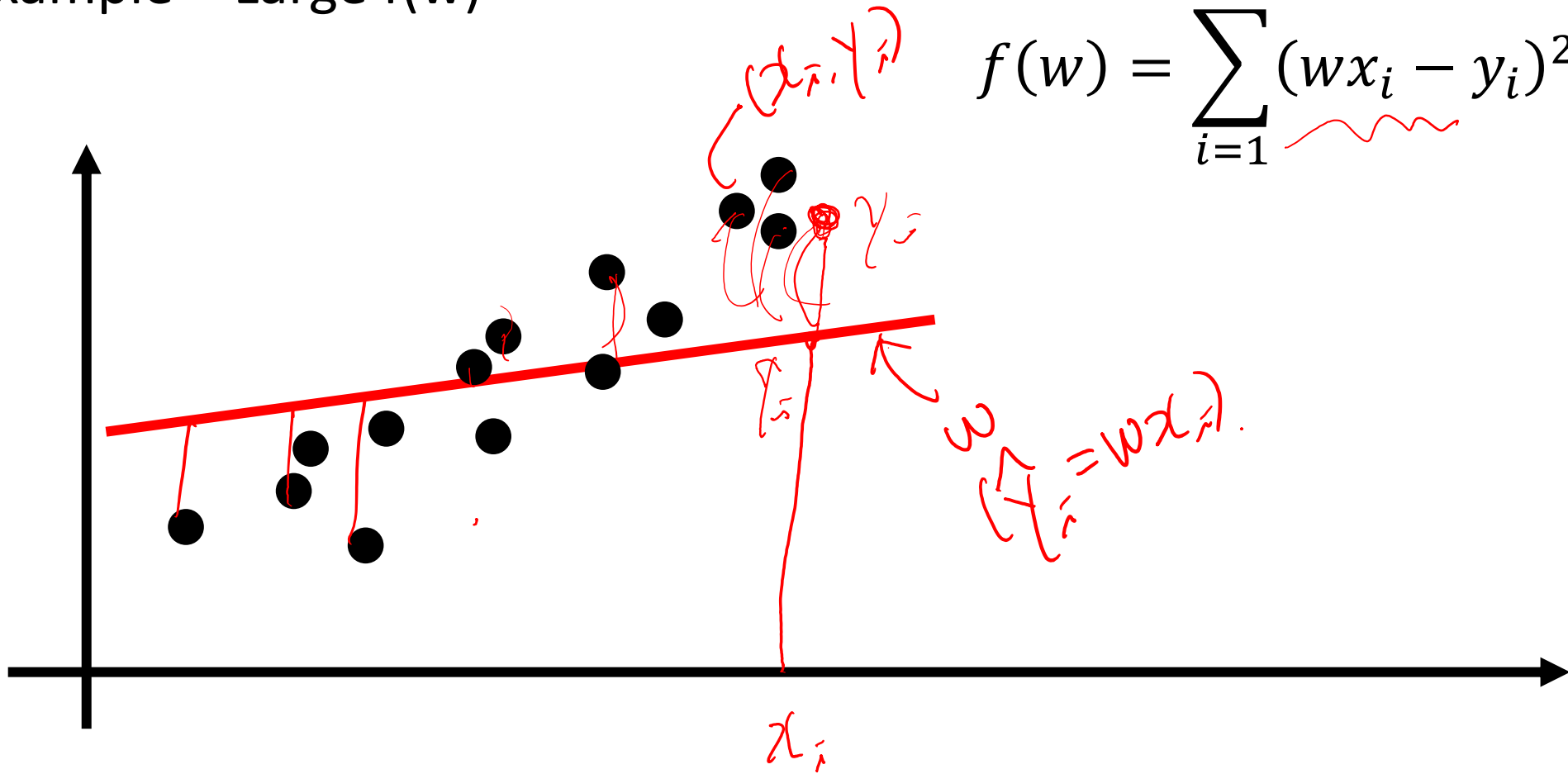- Classic way is setting slope 'w' to minimize the sum of squared errors:

$$\text{avg} \quad \min_{w} f(w) = \sum_{i=1}^{n} (\overset{\hat{y_i}}{wx_i} - y_i)^2$$

  - A probabilistic interpretation is coming later in this course!
  - But usually, it is done because it is easy to minimize.
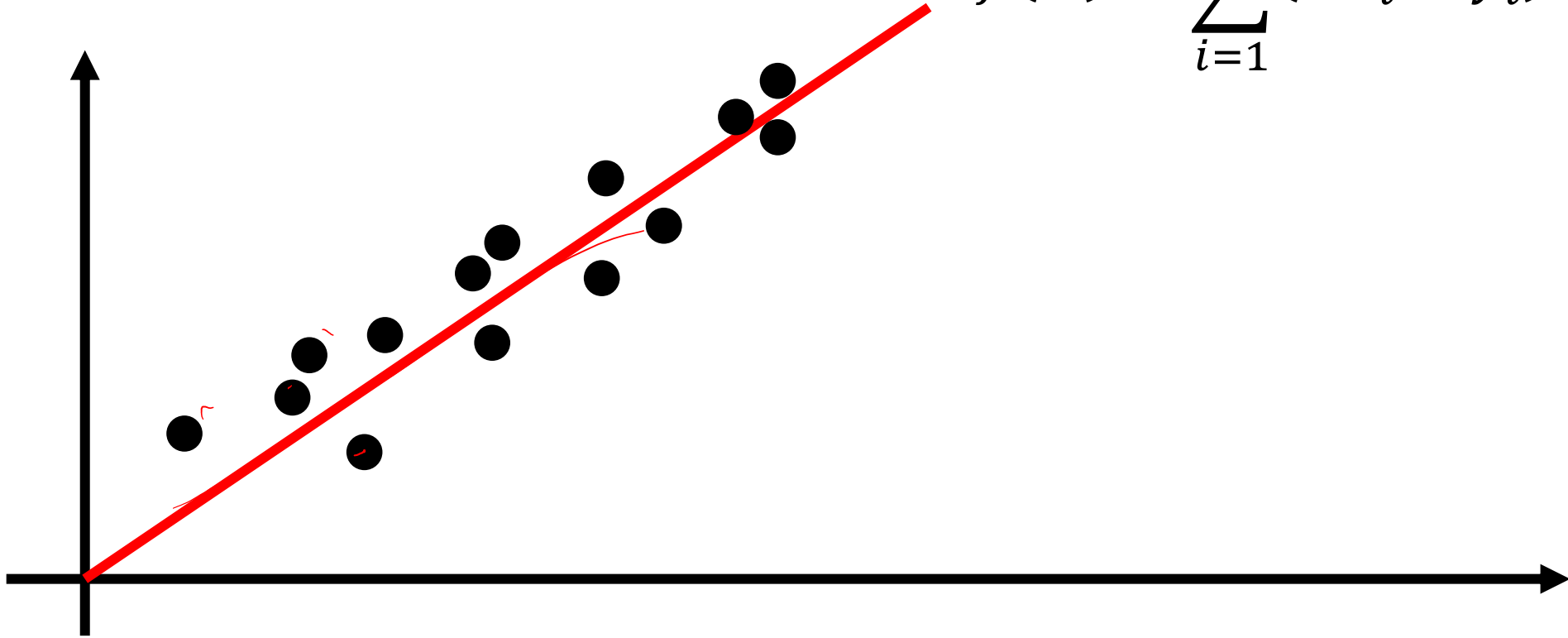
# Least Squres Objective

- Example – Large f(w)

$$f(w) = \sum_{i=1}^{n} (wx_i - y_i)^2$$

# Least Squres Objective

● Example – Small f(w)

$$f(w) = \sum_{i=1}^{n}(wx_i - y_i)^2$$

# Finding Least Squares Solution

- Not change the solution!
  - Multiply 'f' by any positive constant
  - Add some constants to 'f'

$$f'(w) = C_1 \sum_{i=1}^{n} (wx_i - y_i)^2 + C_2$$

- Finding 'w' that minimizes sum of squared errors:

$$f(w) = \frac{1}{2}\sum_{i=1}^{n}(wx_i - y_i)^2 = \frac{1}{2}\sum_{i=1}^{n}\left[w^2 x_i^2 - 2wx_i y_i + y_i^2\right]$$

$$= \frac{w^2}{2}\sum_{i=1}^{n} x_i^2 - w\sum_{i=1}^{n} x_i y_i + \frac{1}{2}\sum_{i=1}^{n} y_i^2 = \frac{w^2}{2}a - wb + c$$

$$f'(w) = wa - b = 0 \qquad w = \frac{\sum_{i=1}^{n} x_i y_i}{\sum_{i=1}^{n} x_i^2}$$

# Multiple Dimension Linear Function

- A simple way is with a d-dimensional linear model
  - $\widehat{y}_i = w_1 x_{i1} + w_2 x_{i2} + w_3 x_{i3} + \cdots + x_d x_{id}$
  - In words, our model is that the output is a weighted sum of the inputs

- We can re-write this in summation notation:
  - $\widehat{y}_i = \sum_{j=1}^{d} w_j x_{ij}$

- We can also re-write this in vector notation: (inner product)
  - $\widehat{y}_i = \mathbf{w}^{\mathrm{T}} \mathbf{x}_i$
  - In this course, a vector is a column vector

# Least Squares in d-Dimensions

- The linear least squares model in d-dimensions minimizes:

$$f(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^{n} (\mathbf{w}^T \mathbf{x}_i - y_i)^2$$
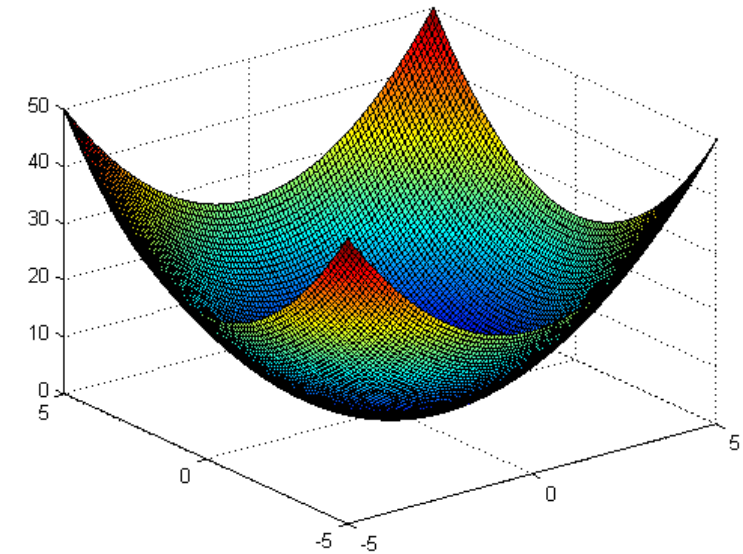
- Least Squares Partial Derivatives for 1 sample

$$f(w_1, w_2, \ldots, w_d) = \frac{1}{2} \left( \sum_{j=1}^{d} w_j x_{ij} \right)^2 - \left( \sum_{j=1}^{d} w_j x_{ij} \right) y_i + \frac{1}{2} y_i^2$$

$$\frac{\partial}{\partial w_k} f(w_1, w_2, \ldots, w_d) = \left( \sum_{j=1}^{d} w_j x_{ij} \right) x_{ik} - y_i x_{ik} = (\mathbf{w}^T \mathbf{x}_i - y_i) x_{ik}$$

$= 0$

# Least Squares in d-Dimensions



- Least Squares Partial Derivatives for all samples

$$\frac{\partial}{\partial w_k} f(w_1, w_2, \ldots, w_d) = \sum_{i=1}^{n} (\mathbf{w}^{\mathrm{T}} \mathbf{x}_i - y_i) x_{ik}$$

- Unfortunately, the partial derivative for $w_j$ depends on all $\{w_1, w_2, \ldots, w_d\}$
  - Thus, we can't just set equal to 0 and solve for $w_j$
  - **We need to find 'w' where the gradient vector equals the zero vector!**

$$\nabla f(w_1, w_2, \ldots, w_d) = \left[ \frac{\partial f}{\partial w_1}, \frac{\partial f}{\partial w_2}, \frac{\partial f}{\partial w_3}, \ldots, \frac{\partial f}{\partial w_d} \right]^{T} = 0$$

# Linear and Quadratic Gradients

$$f(w) = aw^2 \implies \frac{\partial}{\partial w} f(w) = 2aw$$

$$g(w) = bw \implies \frac{\partial}{\partial w} g(w) = b$$

$$h(w) = c \implies \frac{\partial}{\partial w} h(w) = 0$$

$$A = A^T.$$

$$f(\mathbf{w}) = \mathbf{w}^T \mathbf{A} \mathbf{w} \implies \nabla f(\mathbf{w}) = \mathbf{A} \mathbf{w}$$

**If A is symmetric**

$$g(\mathbf{w}) = \mathbf{w}^T \mathbf{b} \implies \nabla g(\mathbf{w}) = \mathbf{b}$$

$$h(\mathbf{w}) = \mathbf{c} \implies \nabla h(\mathbf{w}) = 0$$

# Linear and Quadratic Gradients

- We can re-write the d-dimensional quadratic:
  - $f(\mathbf{w}) = \frac{1}{2}\sum_{i=1}^{n}(\mathbf{w}^T\mathbf{x}_i - y_i)^2 = \frac{1}{2}\|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = \frac{1}{2}\mathbf{w}^T\mathbf{X}^T\mathbf{X}\mathbf{w} - \mathbf{w}^T\mathbf{X}^T\mathbf{y} + \frac{1}{2}\mathbf{y}^T\mathbf{y}$
  - $f(\mathbf{w}) = \frac{1}{2}\mathbf{w}^T\mathbf{A}\mathbf{w} + \mathbf{w}^T\mathbf{b} + \mathbf{c}$

- Thus, the gradient is given by:
  - $\nabla f(\mathbf{w}) = \mathbf{A}\mathbf{w} - \mathbf{b} = \mathbf{X}^T\mathbf{X}\mathbf{w} - \mathbf{X}^T\mathbf{y}$

- **Normal equations:**
  - $\mathbf{X}^T\mathbf{X}\mathbf{w} = \mathbf{X}^T\mathbf{y}$      (in the form of $\mathbf{A}\mathbf{x} = \mathbf{b}$)
  - When $\mathrm{X}^T\mathrm{X}$ is invertible, $\mathbf{w} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$

# Inverse Matrix & Pseudo-inverse Matrix

- For $m \times n$ system $A$,
  - $Ax = b$
  - Let's find $x$ that minimizes the energy of $\|Ax - b\|^2$

- The derivative of $\|Ax - b\|^2$ becomes
  - $Ax - b = 0$
  - $Ax = b$
  - But, we can't estimate the inverse of $A$ because $A$ is not square!

$$w = (X^T X)^{-1} X^T y$$

$100 \times 100$

$X \in \mathbb{R}^{n \times d}$

# Inverse Matrix & Pseudo-inverse Matrix

- The derivative of $\|Ax - b\|^2$ becomes
  - $Ax - b = 0$
  - $Ax = b$
  - But, we can't estimate the inverse of $A$ because $A$ is not square!

- Then, let's make it square matrix
  - $A^T Ax = A^T b$
  - When the columns of A are linearly independent, $A^T A$ is invertible. Thus,
  - $x = (A^T A)^{-1} A^T b \equiv A^+ b$