**CS 1571 Introduction to AI**
**Lecture 24**

- **Density estimation**
- **Linear regression**

**Milos Hauskrecht**
milos@cs.pitt.edu
5329 Sennott Square
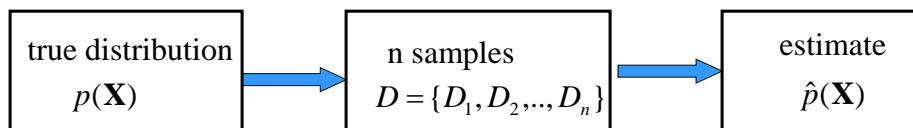
---

# Density estimation

**Data:** $D = \{D_1, D_2, ..., D_n\}$
$D_i = \mathbf{x}_i$      a vector of attribute values

**Objective:** try to estimate the underlying true probability
distribution over variables $\mathbf{X}$, $p(\mathbf{X})$, using examples in $D$

| true distribution $p(\mathbf{X})$ | → | n samples $D = \{D_1, D_2, ..., D_n\}$ | → | estimate $\hat{p}(\mathbf{X})$ |
|---|---|---|---|---|

**Standard (iid) assumptions: Samples**
- **are independent of each other**
- **come from the same (identical) distribution (fixed $p(\mathbf{X})$)**

# Learning via parameter estimation

In this lecture we consider **parametric density estimation**
**Basic settings:**
- A set of random variables $\mathbf{X} = \{X_1, X_2, \ldots, X_d\}$
- **A model of the distribution** over variables in $X$
  with parameters $\Theta$
- **Data** $D = \{D_1, D_2, \ldots, D_n\}$

**Objective:** find parameters $\hat{\Theta}$ that fit the data the best

- What is the best set of parameters?
  – There are various criteria one can apply here.

# Parameter estimation. Basic criteria.

- **Maximum likelihood (ML)**

  maximize $p(D \,|\, \Theta, \xi)$

  $\xi$ - represents prior (background) knowledge

- **Maximum a posteriori probability (MAP)**

  maximize $p(\Theta \,|\, D, \xi)$

  **Selects the mode of the posterior**

  $$p(\Theta \,|\, D, \xi) = \frac{p(D \,|\, \Theta, \xi)\, p(\Theta \,|\, \xi)}{p(D \,|\, \xi)}$$

# Parameter estimation. Coin example.

**Coin example:** we have a coin that can be biased

**Outcomes:** two possible values -- head or tail

**Data:** $D$ a sequence of outcomes $x_i$ such that

- **head** $x_i = 1$
- **tail** $x_i = 0$

**Model:** probability of a head $\theta$

probability of a tail $(1-\theta)$

**Objective:**

We would like to estimate the probability of a **head** $\hat{\theta}$

from data

---

# Parameter estimation. Example.

- **Assume** the unknown and possibly biased coin
- Probability of the head is $\theta$
- **Data:**

  H H T T H H T H T H T T T H T H H H H T H H H H T

  – **Heads:** 15
  – **Tails:** 10

What would be your estimate of the probability of a head ?

$$\tilde{\theta} = ?$$

# Parameter estimation.  Example

- **Assume** the unknown and possibly biased coin
- Probability of the head is $\theta$
- **Data:**

  H H T T H H T H T H T T T H T H H H H T H H H H T
  - **Heads:** 15
  - **Tails:** 10

What would be your choice of the probability of a head ?

**Solution:**  use frequencies of occurrences to do the estimate

$$\tilde{\theta} = \frac{15}{25} = 0.6$$

This is **the maximum likelihood estimate** of the parameter $\theta$

---

# Probability of an outcome

**Data:** $D$   a sequence of outcomes $x_i$  such that
- **head**   $x_i = 1$
- **tail**   $x_i = 0$

**Model:**  probability of a head   $\theta$
probability of a tail   $(1-\theta)$

**Assume: we know the probability** $\theta$
**Probability of an outcome of a coin flip** $x_i$

$$P(x_i \mid \theta) = \theta^{x_i}(1-\theta)^{(1-x_i)} \longleftarrow \quad \textbf{Bernoulli distribution}$$

- Combines the probability of a head and a tail
- So that $x_i$  is going to pick its correct probability
- Gives $\theta$      for $x_i = 1$
- Gives $(1-\theta)$  for $x_i = 0$

# Probability of a sequence of outcomes.

**Data:** $D$   a sequence of outcomes   $x_i$ such that
- **head**   $x_i = 1$
- **tail**   $x_i = 0$

**Model:** probability of a head   $\theta$
     probability of a tail   $(1-\theta)$

**Assume: a sequence of independent coin flips**

     **D = H H T H T H**      **(encoded as D= 110101)**

What is the probability of observing the data sequence **D:**

$$P(D \mid \theta) = ?$$

---

# Probability of a sequence of outcomes.

**Data:** $D$   a sequence of outcomes   $x_i$ such that
- **head**   $x_i = 1$
- **tail**   $x_i = 0$

**Model:** probability of a head   $\theta$
     probability of a tail   $(1-\theta)$

**Assume: a sequence of coin flips D = H H T H T H**

   **encoded as D= 110101**

What is the probability of observing a data sequence **D:**

$$P(D \mid \theta) = \theta\theta\,(1-\theta)\theta\,(1-\theta)\theta$$

## Probability of a sequence of outcomes.

**Data:** $D$ a sequence of outcomes $x_i$ such that
- **head** $x_i = 1$
- **tail** $x_i = 0$

**Model:** probability of a head $\theta$
probability of a tail $(1 - \theta)$

**Assume: a sequence of coin flips D = H H T H T H**

**encoded as D= 110101**

What is the probability of observing a data sequence **D:**

$$P(D \mid \theta) = \theta\theta(1 - \theta)\theta(1 - \theta)\theta$$

**likelihood of the data**

---

## Probability of a sequence of outcomes.

**Data:** $D$ a sequence of outcomes $x_i$ such that
- **head** $x_i = 1$
- **tail** $x_i = 0$

**Model:** probability of a head $\theta$
probability of a tail $(1 - \theta)$

**Assume: a sequence of coin flips D = H H T H T H**

**encoded as D= 110101**

What is the probability of observing a data sequence **D:**

$$P(D \mid \theta) = \theta\theta(1 - \theta)\theta(1 - \theta)\theta$$

$$P(D \mid \theta) = \prod_{i=1}^{6} \theta^{x_i}(1 - \theta)^{(1-x_i)}$$

Can be rewritten using the Bernoulli distribution:

## The goodness of fit to the data.

**Learning: we do not know the value of the parameter** $\theta$

**Our learning goal**:

- Find the parameter $\theta$ that fits the data D the best?

**One solution to the "best":** Maximize the likelihood

$$P(D \mid \theta) = \prod_{i=1}^{n} \theta^{x_i} (1 - \theta)^{(1 - x_i)}$$

**Intuition:**

- more likely are the data given the model, the better is the fit

**Note:** Instead of an error function that measures how bad the data fit the model we have a measure that tells us how well the data fit :

$$Error \ (D, \theta) = -P(D \mid \theta)$$

---

## Maximum likelihood (ML) estimate.

**Likelihood of data:**
$$P(D \mid \theta, \xi) = \prod_{i=1}^{n} \theta^{x_i} (1 - \theta)^{(1 - x_i)}$$

**Maximum likelihood** estimate

$$\theta_{ML} = \arg \max_{\theta} P(D \mid \theta, \xi)$$

**Optimize log-likelihood (the same as maximizing likelihood)**

$$l(D, \theta) = \log P(D \mid \theta, \xi) = \log \prod_{i=1}^{n} \theta^{x_i} (1 - \theta)^{(1 - x_i)} =$$

$$\sum_{i=1}^{n} x_i \log \theta + (1 - x_i) \log(1 - \theta) = \log \theta \underline{\sum_{i=1}^{n} x_i} + \log(1 - \theta) \underline{\sum_{i=1}^{n} (1 - x_i)}$$

$N_1$ - number of heads seen     $N_2$ - number of tails seen

7

# Maximum likelihood (ML) estimate.

**Optimize log-likelihood**

$$l(D,\theta) = N_1 \log\theta + N_2 \log(1-\theta)$$

**Set derivative to zero**

$$\frac{\partial l(D,\theta)}{\partial \theta} = \frac{N_1}{\theta} - \frac{N_2}{(1-\theta)} = 0$$

**Solving**
$$\theta = \frac{N_1}{N_1 + N_2}$$

**ML Solution:** $\quad \theta_{ML} = \dfrac{N_1}{N} = \dfrac{N_1}{N_1 + N_2}$

---

# Maximum likelihood estimate. Example

- **Assume** the unknown and possibly biased coin
- Probability of the head is $\theta$
- **Data:**

  H H T T H H T H T H T T T H T H H H H T H H H H T
    - **Heads:** 15
    - **Tails:** 10

What is the ML estimate of the probability of a head and a tail?

# Maximum likelihood estimate. Example

- Assume the unknown and possibly biased coin
- Probability of the head is $\theta$
- **Data:**

  H H T T H H T H T H T T T H T H H H H T H H H H T
  - **Heads:** 15
  - **Tails:** 10

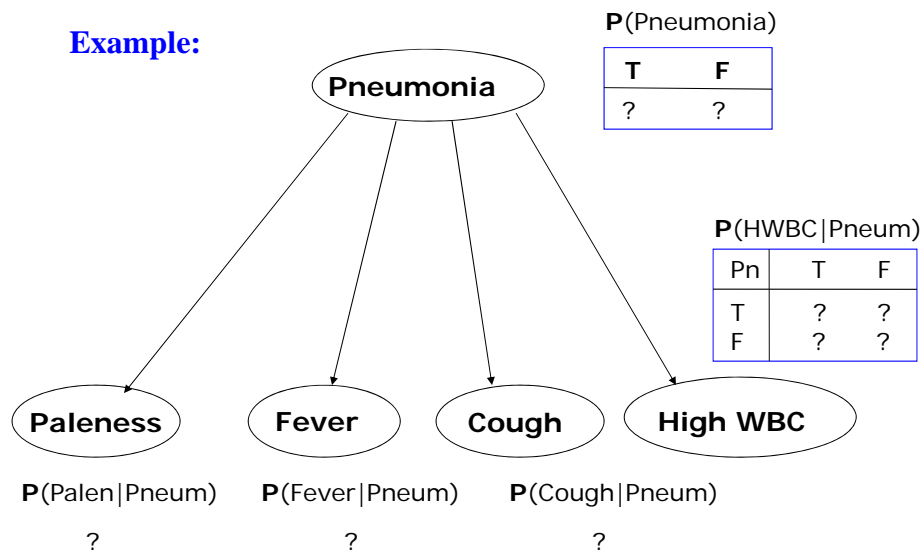What is the ML estimate of the probability of head and tail ?

**Head:** $\quad \theta_{ML} = \dfrac{N_1}{N} = \dfrac{N_1}{N_1 + N_2} = \dfrac{15}{25} = 0.6$

**Tail:** $\quad (1 - \theta_{ML}) = \dfrac{N_2}{N} = \dfrac{N_2}{N_1 + N_2} = \dfrac{10}{25} = 0.4$

---

# Learning of BBN parameters. Example.

**Example:**

**Pneumonia**
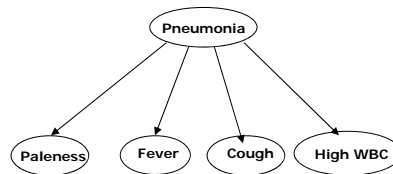
**P**(Pneumonia)

| T | F |
|---|---|
| ? | ? |

**P**(HWBC|Pneum)

| Pn | T | F |
|----|---|---|
| T | ? | ? |
| F | ? | ? |

**Paleness**    **Fever**    **Cough**    **High WBC**

**P**(Palen|Pneum)    **P**(Fever|Pneum)    **P**(Cough|Pneum)

?         ?         ?

## Learning of BBN parameters. Example.

**Data D (different patient cases):**

| Pal | Fev | Cou | HWB | Pneu |
|-----|-----|-----|-----|------|
| T | T | T | T | F |
| T | F | F | F | F |
| F | F | T | T | T |
| F | F | T | F | T |
| F | T | T | T | T |
| T | F | T | F | F |
| F | F | F | F | F |
| T | T | F | F | F |
| T | T | T | T | T |
| F | T | F | T | T |
| T | F | F | T | F |
| F | T | F | F | F |

```
                    Pneumonia
                   /   |    |   \
          Paleness  Fever  Cough  High WBC
```

---

## Estimates of parameters of BBN

- Much like multiple **coin tosses**
- A "smaller" learning problem corresponds to the learning of exactly one conditional distribution
- **Example:**

$$\mathbf{P}(Fever \mid Pneumonia = T)$$
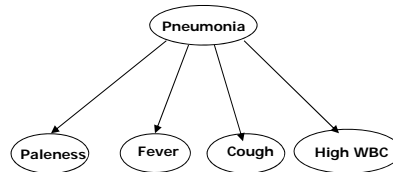
- **Problem:** How to pick the data to learn?

# Learning of BBN parameters. Example.

**Data D (different patient cases):**

| Pal | Fev | Cou | HWB | Pneu |
|-----|-----|-----|-----|------|
| T | T | T | T | F |
| T | F | F | F | F |
| F | F | T | T | T |
| F | F | T | F | T |
| F | T | T | T | T |
| T | F | T | F | F |
| F | F | F | F | F |
| T | T | F | F | F |
| T | T | T | T | T |
| F | T | F | T | T |
| T | F | F | T | F |
| F | T | F | F | F |

Pneumonia → Paleness, Fever, Cough, High WBC

**How to estimate:**

$$\mathbf{P}(Fever \mid Pneumonia = T) = ?$$

---
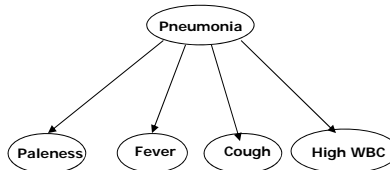
# Learning of BBN parameters. Example.

**Learn:** $\mathbf{P}(Fever \mid Pneumonia = T)$

**Step 1:** Select data points with Pneumonia=T

| Pal | Fev | Cou | HWB | Pneu |
|-----|-----|-----|-----|------|
| T | T | T | T | F |
| T | F | F | F | F |
| F | F | T | T | T |
| F | F | T | F | T |
| F | T | T | T | T |
| T | F | T | F | F |
| F | F | F | F | F |
| T | T | F | F | F |
| T | T | T | T | T |
| F | T | F | T | T |
| T | F | F | T | F |
| F | T | F | F | F |

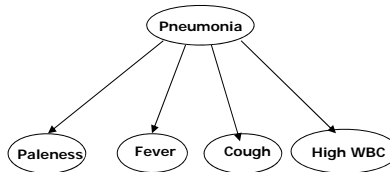Pneumonia → Paleness, Fever, Cough, High WBC

# Learning of BBN parameters. Example.

**Learn:**   $\mathbf{P}(Fever \mid Pneumonia = T)$

**Step 1:**   Ignore the rest

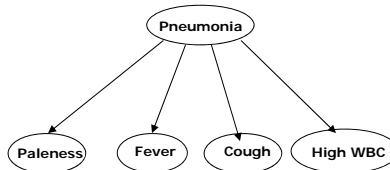| Pal | Fev | Cou | HWB | Pneu |
|-----|-----|-----|-----|------|
| F | F | T | T | T |
| F | F | T | F | T |
| F | T | T | T | T |
| T | T | T | T | T |
| F | T | F | T | T |

---

# Learning of BBN parameters. Example.

**Learn:**   $\mathbf{P}(Fever \mid Pneumonia = T)$

**Step 2:** Select values of the random variable defining the distribution of Fever

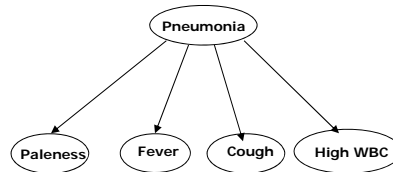| Pal | Fev | Cou | HWB | Pneu |
|-----|-----|-----|-----|------|
| F | F | T | T | T |
| F | F | T | F | T |
| F | T | T | T | T |
| T | T | T | T | T |
| F | T | F | T | T |

# Learning of BBN parameters. Example.

**Learn:** $\mathbf{P}(Fever \mid Pneumonia = T)$

**Step 2:** Ignore the rest

**Fev**
**F**
**F**
**T**
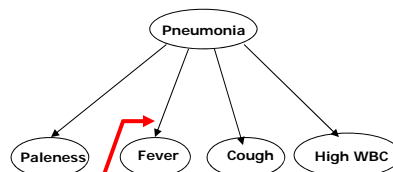**T**
**T**

Pneumonia

Paleness  Fever  Cough  High WBC

---

# Learning of BBN parameters. Example.

**Learn:** $\mathbf{P}(Fever \mid Pneumonia = T)$

**Step 3: Learning the ML estimate**

**Fev**
**F**
**F**
**T**
**T**
**T**

Pneumonia

Paleness  Fever  Cough  High WBC

$\mathbf{P}(Fever \mid Pneumonia = T)$

| T | F |
|-----|-----|
| 0.6 | 0.4 |

# Supervised learning

**Data:** $D = \{D_1, D_2, .., D_n\}$    **a set of *n* examples**

$D_i = <\mathbf{x}_i, y_i>$

$\mathbf{x}_i = (x_{i,1}, x_{i,2}, \cdots x_{i,d})$ is an input vector of size $d$

$y_i$ is the desired output (given by a teacher)

**Objective:** learn the mapping $f : X \rightarrow Y$

s.t. $y_i \approx f(\mathbf{x}_i)$ for all $i = 1, .., n$

- **Regression:** Y is **continuous**
  Example: earnings, product orders $\rightarrow$ company stock price
- **Classification:** Y is **discrete**
  Example: handwritten digit in binary form $\rightarrow$ digit label

---

# Supervised learning

**Next:**

Two basic models of $f : X \rightarrow Y$ used in supervised learning

- **Linear regression:**
  – **Regression where Y is in $\mathcal{R}$**
- **Logistic regression**
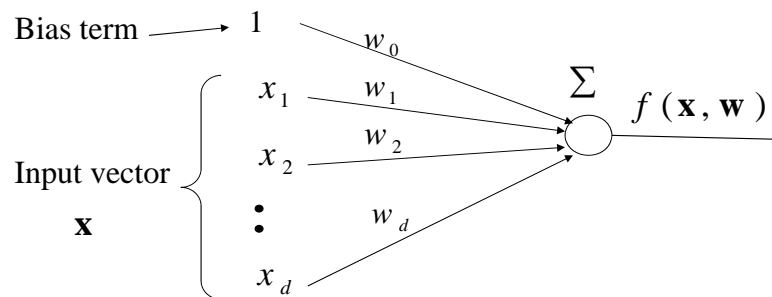  – **Classification with 2 classes**

# Linear regression

- **Function** $f : X \to Y$ is a linear combination of input components

$$f(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + \ldots w_d x_d = w_0 + \sum_{j=1}^{d} w_j x_j$$

$w_0, w_1, \ldots w_k$ - **parameters (weights)**

Bias term $\longrightarrow$ 1   $w_0$

$x_1$   $w_1$

Input vector   $x_2$   $w_2$

$\mathbf{x}$   $\bullet$   $w_d$
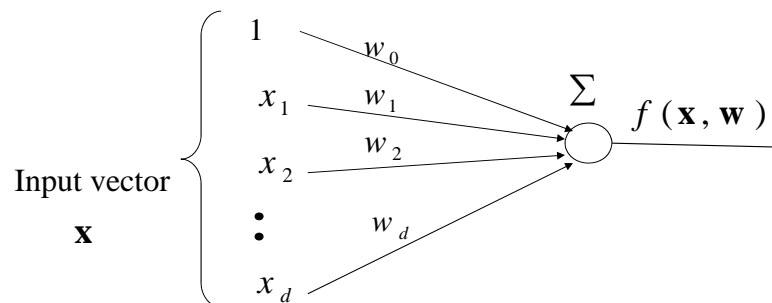
$x_d$

$\Sigma$   $f(\mathbf{x}, \mathbf{w})$

---

# Linear regression

- **Shorter (vector) definition of the model**
  - Include bias constant in the input vector

$$\mathbf{x} = (1, x_1, x_2, \cdots x_d)$$

$$f(\mathbf{x}) = w_0 x_0 + w_1 x_1 + w_2 x_2 + \ldots w_d x_d = \mathbf{w}^T \mathbf{x}$$

$w_0, w_1, \ldots w_k$ - **parameters (weights)**

1   $w_0$

$x_1$   $w_1$

$x_2$   $w_2$

Input vector   $\bullet$   $w_d$

$\mathbf{x}$

$x_d$

$\Sigma$   $f(\mathbf{x}, \mathbf{w})$

# Linear regression. Error.

- **Data:** $D_i = <\mathbf{x}_i, y_i>$
- **Function:** $\mathbf{x}_i \rightarrow f(\mathbf{x}_i)$
- We would like to have $y_i \approx f(\mathbf{x}_i)$ for all $i = 1,.., n$

- **Error function** measures how much our predictions deviate from the desired answers

  **Mean-squared error** $\quad J_n = \dfrac{1}{n} \sum_{i=1,..n} (y_i - f(\mathbf{x}_i))^2$
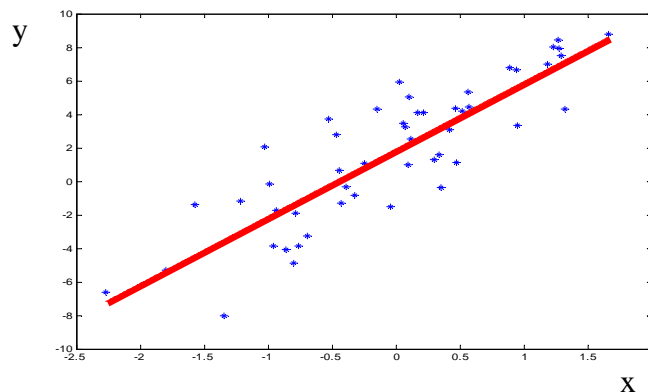
- **Learning:**
    **We want to find the weights minimizing the error !**

---

# Linear regression. Example
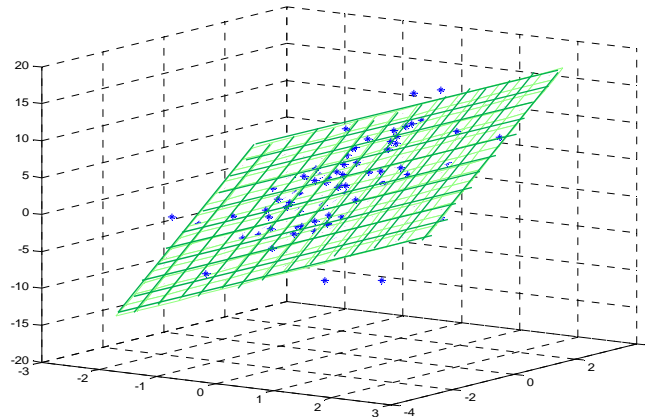
- 1 dimensional input $\qquad \mathbf{x} = (x_1)$

# Linear regression. Example.

- 2 dimensional input $\quad \mathbf{x} = (x_1, x_2)$

# Linear regression. Optimization.

- We want the **weights minimizing the error**

$$J_n = \frac{1}{n} \sum_{i=1,..n} (y_i - f(\mathbf{x}_i))^2 = \frac{1}{n} \sum_{i=1,..n} (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

- For the optimal set of parameters, derivatives of the error with respect to each parameter must be 0

$$\frac{\partial}{\partial w_j} J_n(\mathbf{w}) = -\frac{2}{n} \sum_{i=1}^{n} (y_i - w_0 x_{i,0} - w_1 x_{i,1} - \ldots - w_d x_{i,d}) x_{i,j} = 0$$

- **Vector of derivatives:**

$$\text{grad}_{\mathbf{w}}(J_n(\mathbf{w})) = \nabla_{\mathbf{w}}(J_n(\mathbf{w})) = -\frac{2}{n} \sum_{i=1}^{n} (y_i - \mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i = \overline{\mathbf{0}}$$

## Linear regression. Optimization.

- For the optimal set of parameters, derivatives of the error with respect to each parameter must be 0

$$J_n = \frac{1}{n}\sum_{i=1,..n}(y_i - f(\mathbf{x}_i))^2 = \frac{1}{n}\sum_{i=1,..n}(y_i - [w_0 + w_1 x^{(1)} + w_2 x^{(2)} + \ldots w_k x^{(k)}])^2$$

- $\mathrm{grad}_{\,\mathbf{w}}(J_n(\mathbf{w})) = \overline{\mathbf{0}}$  defines a set of equations in $\mathbf{w}$

$$\frac{\partial}{\partial w_0} J_n(w) = -\frac{2}{n}\sum_{i=1}^{n}[y_i - (w_0 + w_1 x^{(1)} + w_2 x^{(2)} + \ldots w_k x^{(k)})] = 0$$

$$\frac{\partial}{\partial w_1} J_n(w) = -\frac{2}{n}\sum_{i=1}^{n}[y_i - (w_0 + w_1 x^{(1)} + w_2 x^{(2)} + \ldots w_k x^{(k)})]x^{(1)} = 0$$

$$\ldots$$

$$\frac{\partial}{\partial w_j} J_n(w) = -\frac{2}{n}\sum_{i=1}^{n}[y_i - (w_0 + w_1 x^{(1)} + w_2 x^{(2)} + \ldots w_k x^{(k)})]x^{(j)} = 0$$

$$\ldots$$

---

## Solving linear regression

$$\frac{\partial}{\partial w_j} J_n(\mathbf{w}) = -\frac{2}{n}\sum_{i=1}^{n}(y_i - w_0 x_{i,0} - w_1 x_{i,1} - \ldots - w_d x_{i,d})x_{i,j} = 0$$

By rearranging the terms we get a **system of linear equations** with $d+1$ unknowns

$$\mathbf{A}\mathbf{w} = \mathbf{b}$$

$$w_0\sum_{i=1}^{n}x_{i,0}1 + w_1\sum_{i=1}^{n}x_{i,1}1 + \ldots + w_j\sum_{i=1}^{n}x_{i,j}1 + \ldots + w_d\sum_{i=1}^{n}x_{i,d}1 = \sum_{i=1}^{n}y_i 1$$

$$w_0\sum_{i=1}^{n}x_{i,0}x_{i,1} + w_1\sum_{i=1}^{n}x_{i,1}x_{i,1} + \ldots + w_j\sum_{i=1}^{n}x_{i,j}x_{i,1} + \ldots + w_d\sum_{i=1}^{n}x_{i,d}x_{i,1} = \sum_{i=1}^{n}y_i x_{i,1}$$

$$\bullet\bullet\bullet$$

$$w_0\sum_{i=1}^{n}x_{i,0}x_{i,j} + w_1\sum_{i=1}^{n}x_{i,1}x_{i,j} + \ldots + w_j\sum_{i=1}^{n}x_{i,j}x_{i,j} + \ldots + w_d\sum_{i=1}^{n}x_{i,d}x_{i,j} = \sum_{i=1}^{n}y_i x_{i,j}$$

$$\bullet\bullet\bullet$$

# Solving linear regression

- The optimal set of weights satisfies:

$$\nabla_{\mathbf{w}}(J_n(\mathbf{w})) = -\frac{2}{n}\sum_{i=1}^{n}(y_i - \mathbf{w}^T\mathbf{x}_i)\mathbf{x}_i = \overline{\mathbf{0}}$$

Leads to a **system of linear equations (SLE)** with $d+1$
unknowns of the form $\mathbf{Aw} = \mathbf{b}$

$$w_0\sum_{i=1}^{n}x_{i,0}x_{i,j} + w_1\sum_{i=1}^{n}x_{i,1}x_{i,j} + \ldots + w_j\sum_{i=1}^{n}x_{i,j}x_{i,j} + \ldots + w_d\sum_{i=1}^{n}x_{i,d}x_{i,j} = \sum_{i=1}^{n}y_ix_{i,j}$$

**Solutions to SLE:**
- e.g. matrix inversion (if the matrix is singular)

$$\mathbf{w} = \mathbf{A}^{-1}\mathbf{b}$$

---

# Gradient descent solution

- There are other ways to solve the weight optimization problem in the linear regression model

$$J_n = Error(\mathbf{w}) = \frac{1}{n}\sum_{i=1,\ldots n}(y_i - f(\mathbf{x}_i, \mathbf{w}))^2$$

- A simple technique:
  – **Gradient descent**

  **Idea:**
  - Adjust weights in the direction that improves the Error
  - The gradient tells us what is the right direction

  $$\mathbf{w} \leftarrow \mathbf{w} - \alpha\,\nabla_{\mathbf{w}}Error_i(\mathbf{w})$$

  $\alpha > 0$ - a learning rate (scales the gradient changes)