# Foundations of Machine Learning
## Linear Regression (contd)

August 25, 2020

## Recap: Least Squares Linear Regression

Given: Training data containing $n$ samples where each $\mathbf{x}_i \in \mathbb{R}^d$, $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^{n}$

Least squares solution to linear regression is as follows:

$$\mathbf{w}^* = \arg\min_{\mathbf{w}} E(\mathbf{w}, \mathcal{D}) = \arg\min_{\mathbf{w}} \left\{ \sum_{i=1}^{n} \left( \sum_{j=0}^{m} w_j \phi_j(\mathbf{x}_i) - y_i \right)^2 \right\}$$

$$\mathbf{w}^* = \arg\min_{\mathbf{w}} \left( ||\Phi\mathbf{w} - \mathbf{y}||^2 \right) = (\Phi^T\Phi)^{-1}\Phi^T\mathbf{y}$$

where

$$\mathbf{\Phi} = \begin{bmatrix} 1 & \phi_1(\mathbf{x}_1) & \dots & \phi_m(\mathbf{x}_1) \\ \dots & \dots & \dots \\ 1 & \phi_1(\mathbf{x}_n) & \dots & \phi_m(\mathbf{x}_n) \end{bmatrix}, \mathbf{w} = \begin{bmatrix} w_0 \\ \dots \\ w_m \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y_1 \\ \dots \\ y_n \end{bmatrix}$$

$n \times (m+1)$     $(m+1) \times 1$     $n \times 1$

## Probabilistic Modeling of Linear Regression

- From the training set $\mathcal{D} = \{\mathbf{x}_j, y_j\}_{j=1}^n$, assume each target $y_j$ is drawn from a probabilistic model.

## Probabilistic Modeling of Linear Regression

- From the training set $\mathcal{D} = \{\mathbf{x}_j, y_j\}_{j=1}^n$, assume each target $y_j$ is drawn from a probabilistic model.

- Let us assume a noisy linear model parameterized by $\mathbf{w} \in \mathbb{R}^d$

$$y_j = \mathbf{w}^T \mathbf{x}_j + \epsilon_j$$

  Each target is associated with an additive stochastic noise modeled as a random variable $\epsilon_j$.

## Probabilistic Modeling of Linear Regression

- From the training set $\mathcal{D} = \{\mathbf{x}_j, y_j\}_{j=1}^n$, assume each target $y_j$ is drawn from a probabilistic model.

- Let us assume a noisy linear model parameterized by $\mathbf{w} \in \mathbb{R}^d$

$$y_j = \mathbf{w}^T \mathbf{x}_j + \epsilon_j$$

Each target is associated with an additive stochastic noise modeled as a random variable $\epsilon_j$.

- Let us model the noise to come from a **Gaussian distribution:**
  $\epsilon_j \sim \mathcal{N}(\underset{\text{mean}}{\underline{0}}, \underset{}{\underline{\sigma^2}})$. variance

3

## Probabilistic Modeling of Linear Regression

- From the training set $\mathcal{D} = \{\mathbf{x}_j, y_j\}_{j=1}^{n}$, assume each target $y_j$ is drawn from a probabilistic model.

- Let us assume a noisy linear model parameterized by $\mathbf{w} \in \mathbb{R}^d$

$$y_j = \mathbf{w}^T \mathbf{x}_j + \epsilon_j$$

  Each target is associated with an additive stochastic noise modeled as a random variable $\epsilon_j$.

- Let us model the noise to come from a **Gaussian distribution:** $\epsilon_j \sim \mathcal{N}(0, \sigma^2)$.

- Noise variables $\epsilon_j$ all have the same mean (0), same variance ($\sigma^2$) and $\mathrm{Cov}(\epsilon_i, \epsilon_j) = 0$ for $i \neq j$.

## Probabilistic Modeling of Linear Regression

- If $y_j = \mathbf{w}^T\mathbf{x}_j + \epsilon_j$ and $\epsilon_j \sim \mathcal{N}(0, \sigma^2)$, then each $y_j$ is drawn from the following Gaussian:

$$y_j \sim \mathcal{N}(\mathbf{w}^T\mathbf{x}_j, \sigma^2)$$

## Probabilistic Modeling of Linear Regression

- If $y_j = \mathbf{w}^T\mathbf{x}_j + \epsilon_j$ and $\epsilon_j \sim \mathcal{N}(0, \sigma^2)$, then each $y_j$ is drawn from the following Gaussian:

$$y_j \sim \mathcal{N}(\mathbf{w}^T\mathbf{x}_j, \sigma^2)$$

- Objective: Estimate $\mathbf{w}$ from this probabilistic model. We will look at **Maximum Likelihood Estimation (MLE)** in this lecture.

## MLE, more generally

A set of independent and identically distributed (i.i.d.) observations
$\mathbf{y} = \{y_1, \ldots, y_n\}$ are generated according to a probability model parameterized by
$\theta$, i.e. $y_j \sim \Pr(\mathbf{y}|\theta)$

## MLE, more generally

A set of independent and identically distributed (i.i.d.) observations
$\mathbf{y} = \{y_1, \ldots, y_n\}$ are generated according to a probability model parameterized by
$\theta$, i.e. $y_j \sim \Pr(\mathbf{y}|\theta)$

What is $\Pr(\mathbf{y}|\theta)$ (also referred to as the *likelihood*)?

## MLE, more generally

A set of independent and identically distributed (i.i.d.) observations
$\mathbf{y} = \{y_1, \ldots, y_n\}$ are generated according to a probability model parameterized by
$\theta$, i.e. $y_j \sim \Pr(\mathbf{y}|\theta)$

What is $\Pr(\mathbf{y}|\theta)$ (also referred to as the *likelihood*)?

Find the parameters $\theta$. What's one good way? Find $\theta$ s.t. it maximizes the
likelihood of the observed data. Hence, **MLE**.

## MLE contd

- Typically maximize log-likelihood i.e. $\log(\Pr(\mathbf{y}|\theta))$ instead of the likelihood. Would this affect the estimated parameters?

## MLE contd

- Typically maximize log-likelihood i.e. $\log(\Pr(\mathbf{y}|\theta))$ instead of the likelihood. Would this affect the estimated parameters?

- Why apply this log transformation? Mathematical convenience
  Numerical convenience

## MLE contd

- Typically maximize log-likelihood i.e. $\log(\Pr(\mathbf{y}|\theta))$ instead of the likelihood. Would this affect the estimated parameters?

- Why apply this log transformation?

$$\Pr(y_j ; \theta) \quad \text{or} \quad P_\theta(y_j)$$

- What is the log-likelihood?

$$\mathcal{L}(\theta) = \log(\Pr(\mathbf{y}|\theta)) = \log(\prod_{j=1}^{n} \Pr(y_j|\theta)) = \sum_{j=1}^{n} \log(\Pr(y_j|\theta))$$

- The MLE estimate is given by:

$$\theta_{\mathsf{MLE}} = \underset{\theta}{\mathrm{argmax}}\, \mathcal{L}(\theta) = \underset{\theta}{\mathrm{argmax}} \sum_{j=1}^{n} \log(\Pr(y_j|\theta))$$

# Illustrating MLE using a simple coin toss example

Say we toss a coin $N$ times. Each coin toss $y_j$ is a binary random variable with a Bernoulli distribution $\Pr(y_j|\theta) = \underline{\theta}^{y_j}(1-\theta)^{1-y_j}$. Estimate $\theta$ using MLE.

$$L(\theta) = \log \prod_{j=1}^{N} P(y_j|\theta) = \sum_j \log P(y_j|\theta)$$

$$= \sum_j y_j \log\theta + (1-y_j)\log(1-\theta)$$

$$\frac{\partial L(\theta)}{\partial\theta} = 0 \implies \sum_j \frac{y_j}{\theta} - \frac{(1-y_j)}{1-\theta} = 0 \implies \boxed{\theta = \frac{\sum_j y_j}{N}}$$

## Coming back to MLE for linear regression

Assuming that each target $y_j$ is Gaussian-distributed, i.e.

$$\Pr(y_j|x_j, \mathbf{w}) = \mathcal{N}(\mathbf{w}^T x_j, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{ \frac{-(y_j - \mathbf{w}^T x_j)^2}{2\sigma^2} \right\}$$

Likelihood of the data $\mathcal{L}(\mathbf{w}) = \Pr(\mathbf{y}|\mathbf{X}, \mathbf{w})$

Assuming that each target $y_j$ is Gaussian-distributed, i.e.

$$\text{Pr}(y_j|x_j, \mathbf{w}) = \mathcal{N}(\mathbf{w}^T x_j, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{ \frac{-(y_j - \mathbf{w}^T x_j)^2}{2\sigma^2} \right\}$$

Likelihood of the data $\mathcal{L}(\mathbf{w}) = \text{Pr}(\mathbf{y}|\mathbf{X}, \mathbf{w})$

$$\mathcal{L}(w) = \log \prod_{j=1}^{n} P(y_j | x_j, w) = \sum_j \log P(y_j | x_j, w)$$

MLE estimate: $\mathbf{w}_{\text{MLE}} = \underset{\mathbf{w}}{\text{argmax}}\, \mathcal{L}(\mathbf{w})$

## MLE estimate for probabilistic linear regression

To estimate $\mathbf{w}_{\text{MLE}}$, consider log-likelihood $(\log(\mathcal{L}(\mathbf{w})))$ instead of likelihood $(\mathcal{L}(\mathbf{w}))$:

$\log(\mathcal{L}(\mathbf{w})) = LL(\mathbf{w})$

## MLE estimate for probabilistic linear regression

To estimate $\mathbf{w}_{MLE}$, consider log-likelihood ($\log(\mathcal{L}(\mathbf{w}))$) instead of likelihood ($\mathcal{L}(\mathbf{w})$):

$\log(\mathcal{L}(\mathbf{w})) = LL(\mathbf{w})$

$$LL(w) = \sum_j \log P(y_j \mid x_j, w)$$

$$= \text{constant} - \sum_j \frac{(y_j - w^T x_j)^2}{2\sigma^2}$$

$\mathbf{w}_{MLE} = \underset{\mathbf{w}}{\arg\max}\, LL(\mathbf{w}) = \underset{\mathbf{w}}{\arg\min} \sum_{j=1}^{n} (y_j - \mathbf{w}^T x_j)^2$ (**Same as least squares!**)

## Estimating w

Instead of finding the least squares/MLE solution analytically, use a search algorithm that progressively decreases the error function $E(\mathbf{w}, \mathcal{D})$.

- Initialize <u>**w**</u>
- Until <u>we converge</u>:
    Find a new value of **w** that reduces $E(\mathbf{w}, \mathcal{D})$

$$w_{t+1} \Leftarrow w_t + \boxed{\phantom{xxxxxx}}$$

## Estimating w

Instead of finding the least squares/MLE solution analytically, use a search algorithm that progressively decreases the error function $E(\mathbf{w}, \mathcal{D})$.

- Initialize **w**
- Until we converge:
    Find a new value of **w** that reduces $E(\mathbf{w}, \mathcal{D})$

**Gradient descent** is an iterative algorithm used to find the minimum of a function.

# Gradient Descent

$$W = \begin{bmatrix} w_1 \\ \vdots \\ w_t \\ \vdots \\ w_d \end{bmatrix}$$

Initialize $\underline{\mathbf{w}} \leftarrow \vec{0}$ vector

Repeat until convergence:

$w_i \leftarrow w_i - \eta \frac{\partial E(\mathbf{w})}{\partial w_i}$ (Simultaneously update each $w_i$)

$\|\nabla_w E(w)\|^2 \leq \epsilon$

$\|E(W_{t+1}) - E(W_t)\|_2 \leq \epsilon$

Compute $\frac{\partial E(\mathbf{w})}{\partial w_i}$

$$W \leftarrow W - \eta \nabla_w E(w)$$

LEARNING RATE

## Using Gradient Descent (GD)

Weight update rule in GD: $\mathbf{w} \leftarrow \mathbf{w} - \eta\nabla_{\mathbf{w}}E(\mathbf{w})$

$$\nearrow \sum_i E_i(w)$$

Faster version: **Stochastic Gradient Descent (SGD)**

The weight update rule for SGD: $\mathbf{w} \leftarrow \mathbf{w} - \eta\nabla_{\mathbf{w}}E(\mathbf{w}; \mathbf{x}_i, y_i)$ where the loss is computed for a single example randomly sampled from the training set.

Middle ground: **Mini-batch Gradient Descent**

The weight update rule becomes: $\mathbf{w} \leftarrow \mathbf{w} - \eta\nabla_{\mathbf{w}}E_B(\mathbf{w}; \{\mathbf{x}_i, y_i\}_{i=1}^{|B|})$ where the loss is computed for a batch of $|B|$ examples sampled from the training set.