

On Linear Regression: Regularization, Probabilistic Interpretation, and Gradient Descent

Karl Stratos

June 19, 2018

Recall: Linear Regression

- ▶ Given $S = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$, find model parameter $\mathbf{w} \in \mathbb{R}^d$ that minimizes the **sum of squared errors**:

$$J_S^{\text{LS}}(\mathbf{w}) := \sum_{i=1}^n \left(y^{(i)} - \mathbf{w} \cdot \mathbf{x}^{(i)} \right)^2$$

- ▶ We will discuss three topics through linear regression.
 1. **Regularization** to prevent overfitting
 2. **Maximum likelihood estimation** (MLE) interpretation
 3. **Gradient descent** to estimate model parameter
- ▶ Far-reaching implications beyond linear regression

Overview

Regularization

A Probabilistic Interpretation of Linear Regression

Optimization by Local Search

Motivation

- ▶ The least squares solution is the *best* linear regressor on training data S (solved in closed-form):

$$\mathbf{w}_S^{\text{LS}} := \arg \min_{\mathbf{w} \in \mathbb{R}^d} J_S^{\text{LS}}(\mathbf{w})$$

Motivation

- ▶ The least squares solution is the *best* linear regressor on training data S (solved in closed-form):

$$\mathbf{w}_S^{\text{LS}} := \arg \min_{\mathbf{w} \in \mathbb{R}^d} J_S^{\text{LS}}(\mathbf{w})$$

- ▶ But we care nothing about how well we do on S ! Rather, what we really care about is:

Can \mathbf{w}_S^{LS} handle a **new** x not already seen in S ?

Motivation

- ▶ The least squares solution is the *best* linear regressor on training data S (solved in closed-form):

$$\mathbf{w}_S^{\text{LS}} := \arg \min_{\mathbf{w} \in \mathbb{R}^d} J_S^{\text{LS}}(\mathbf{w})$$

- ▶ But we care nothing about how well we do on S ! Rather, what we really care about is:

Can \mathbf{w}_S^{LS} handle a **new** x not already seen in S ?

- ▶ This is the heart of machine learning: theory/applications of *generalization*.

A Model of Noisy Environment

- ▶ There is some “true” parameter $w^* \in \mathbb{R}^d$.

A Model of Noisy Environment

- ▶ There is some “true” parameter $\mathbf{w}^* \in \mathbb{R}^d$.
- ▶ There is some input distribution $\mathbf{x} \sim \mathcal{D}$ and some **noise** distribution $\epsilon \sim \mathcal{E}$. Assume that $\mathbf{E}[\epsilon] = 0$ and $\text{Var}(\epsilon) = \sigma^2$.

A Model of Noisy Environment

- ▶ There is some “true” parameter $\mathbf{w}^* \in \mathbb{R}^d$.
- ▶ There is some input distribution $\mathbf{x} \sim \mathcal{D}$ and some **noise** distribution $\epsilon \sim \mathcal{E}$. Assume that $\mathbf{E}[\epsilon] = 0$ and $\text{Var}(\epsilon) = \sigma^2$.
- ▶ Each sample $(\mathbf{x}, y) \in \mathbb{R}^d \times \mathbb{R}$ is generated by drawing $\mathbf{x} \sim \mathcal{D}$ and $\epsilon \sim \mathcal{E}$ and setting

$$y = \mathbf{w}^* \cdot \mathbf{x} + \epsilon$$

(Thus the training data S is a random variable.)

A Model of Noisy Environment

- ▶ There is some “true” parameter $\mathbf{w}^* \in \mathbb{R}^d$.
- ▶ There is some input distribution $\mathbf{x} \sim \mathcal{D}$ and some **noise** distribution $\epsilon \sim \mathcal{E}$. Assume that $\mathbf{E}[\epsilon] = 0$ and $\text{Var}(\epsilon) = \sigma^2$.
- ▶ Each sample $(\mathbf{x}, y) \in \mathbb{R}^d \times \mathbb{R}$ is generated by drawing $\mathbf{x} \sim \mathcal{D}$ and $\epsilon \sim \mathcal{E}$ and setting

$$y = \mathbf{w}^* \cdot \mathbf{x} + \epsilon$$

(Thus the training data S is a random variable.)

- ▶ Check that \mathbf{w}_S^{LS} is consistent/unbiased:

$$\mathbf{E}_S[\mathbf{w}_S^{\text{LS}}] = \mathbf{E}_S[\mathbf{X}_S^+ (\mathbf{X}_S \mathbf{w}^* + \epsilon)] = \mathbf{w}^*$$

Measuring the Future Performance

- ▶ We want w_S^{LS} which is trained on S to incur small loss in expectation (“true/population error”):

$$\mathbf{E}_{S, \mathbf{x}, \epsilon} \left[\left((\mathbf{w}^* \cdot \mathbf{x} + \epsilon) - w_S^{\text{LS}} \cdot \mathbf{x} \right)^2 \right]$$

Measuring the Future Performance

- ▶ We want w_S^{LS} which is trained on S to incur small loss in expectation (“true/population error”):

$$\mathbf{E}_{S, \mathbf{x}, \epsilon} \left[\left((\mathbf{w}^* \cdot \mathbf{x} + \epsilon) - \mathbf{w}_S^{\text{LS}} \cdot \mathbf{x} \right)^2 \right]$$

- ▶ By the bias-variance decomposition of squared error, this is (omitting the expectation over \mathbf{x}):

$$\underbrace{\left(\mathbf{w}^* \cdot \mathbf{x} - \mathbf{E}_S[\mathbf{w}_S^{\text{LS}}] \cdot \mathbf{x} \right)^2}_{0 \text{ in this case}} + \text{Var}_S \left(\mathbf{w}_S^{\text{LS}} \cdot \mathbf{x} \right) + \underbrace{\sigma^2}_{\text{can't help}}$$

Measuring the Future Performance

- ▶ We want w_S^{LS} which is trained on S to incur small loss in expectation (“true/population error”):

$$\mathbf{E}_{S, \mathbf{x}, \epsilon} \left[\left((w^* \cdot \mathbf{x} + \epsilon) - w_S^{\text{LS}} \cdot \mathbf{x} \right)^2 \right]$$

- ▶ By the bias-variance decomposition of squared error, this is (omitting the expectation over \mathbf{x}):

$$\underbrace{\left(w^* \cdot \mathbf{x} - \mathbf{E}_S[w_S^{\text{LS}}] \cdot \mathbf{x} \right)^2}_{0 \text{ in this case}} + \text{Var}_S(w_S^{\text{LS}} \cdot \mathbf{x}) + \underbrace{\sigma^2}_{\text{can't help}}$$

- ▶ The variance term can be large if parameter values are large.
 - ▶ $(w_S^{\text{LS}} \cdot \mathbf{x})^2$ more sensitive to a perturbation of S

Ridge Regression

- ▶ “Shrink” the size of the estimator by penalizing its l_2 norm:

$$\mathbf{w}_{S,\lambda}^{\text{LSR}} := \arg \min_{\mathbf{w} \in \mathbb{R}^d} J_S^{\text{LS}}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$

Ridge Regression

- ▶ “Shrink” the size of the estimator by penalizing its l_2 norm:

$$\mathbf{w}_{S,\lambda}^{\text{LSR}} := \arg \min_{\mathbf{w} \in \mathbb{R}^d} J_S^{\text{LS}}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$

- ▶ Closed-form solution given by (hence the name)

$$\mathbf{w}_{S,\lambda}^{\text{LSR}} = (\mathbf{X}_S^\top \mathbf{X}_S + \lambda I_{d \times d})^{-1} \mathbf{X}_S^\top \mathbf{y}$$

Ridge Regression

- ▶ “Shrink” the size of the estimator by penalizing its l_2 norm:

$$\mathbf{w}_{S,\lambda}^{\text{LSR}} := \arg \min_{\mathbf{w} \in \mathbb{R}^d} J_S^{\text{LS}}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$

- ▶ Closed-form solution given by (hence the name)

$$\mathbf{w}_{S,\lambda}^{\text{LSR}} = (\mathbf{X}_S^\top \mathbf{X}_S + \lambda I_{d \times d})^{-1} \mathbf{X}_S^\top \mathbf{y}$$

- ▶ No longer unbiased: $\mathbf{E}_S[\mathbf{w}_{S,\lambda}^{\text{LSR}}] \neq \mathbf{w}^*$ for $\lambda > 0$.

Ridge Regression

- ▶ “Shrink” the size of the estimator by penalizing its l_2 norm:

$$\mathbf{w}_{S,\lambda}^{\text{LSR}} := \arg \min_{\mathbf{w} \in \mathbb{R}^d} J_S^{\text{LS}}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$

- ▶ Closed-form solution given by (hence the name)

$$\mathbf{w}_{S,\lambda}^{\text{LSR}} = (\mathbf{X}_S^\top \mathbf{X}_S + \lambda I_{d \times d})^{-1} \mathbf{X}_S^\top \mathbf{y}$$

- ▶ No longer unbiased: $\mathbf{E}_S[\mathbf{w}_{S,\lambda}^{\text{LSR}}] \neq \mathbf{w}^*$ for $\lambda > 0$.

- ▶ But the true error might be smaller!

$$\underbrace{(\mathbf{w}^* \cdot \mathbf{x} - \mathbf{E}_S[\mathbf{w}_{S,\lambda}^{\text{LSR}}] \cdot \mathbf{x})^2}_{\text{no longer 0}} + \underbrace{\text{Var}_S(\mathbf{w}_{S,\lambda}^{\text{LSR}} \cdot \mathbf{x})}_{\text{smaller}} + \underbrace{\sigma^2}_{\text{can't help}}$$

Lasso Regression

- ▶ Another idea: penalize the l_1 norm:

$$\mathbf{w}_{S,\lambda}^{\text{LSL}} := \arg \min_{\mathbf{w} \in \mathbb{R}^d} J_S^{\text{LS}}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

Lasso Regression

- ▶ Another idea: penalize the l_1 norm:

$$\mathbf{w}_{S,\lambda}^{\text{LSL}} := \arg \min_{\mathbf{w} \in \mathbb{R}^d} J_S^{\text{LS}}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

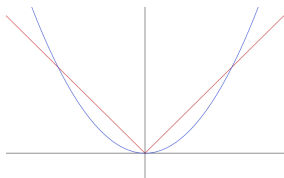
- ▶ Still convex though not differentiable. Can be solved by existing convex optimization methods or subgradient descent.

Lasso Regression

- ▶ Another idea: penalize the l_1 norm:

$$\mathbf{w}_{S,\lambda}^{\text{LSL}} := \arg \min_{\mathbf{w} \in \mathbb{R}^d} J_S^{\text{LS}}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

- ▶ Still convex though not differentiable. Can be solved by existing convex optimization methods or subgradient descent.
- ▶ Solutions with zero entries are encouraged (hence the name).



(squared l_2 norm penalty vs l_1 norm penalty)

Summary on Regularization

- ▶ The l_2/l_1 regularized solutions can be framed as constrained solutions: for some $\alpha, \beta \in \mathbb{R}$

$$\mathbf{w}_{S,\lambda}^{\text{LSL}} := \arg \min_{\mathbf{w} \in \mathbb{R}^d: \|\mathbf{w}\|_2 \leq \alpha} J_S^{\text{LS}}(\mathbf{w})$$

$$\mathbf{w}_{S,\lambda}^{\text{LSR}} := \arg \min_{\mathbf{w} \in \mathbb{R}^d: \|\mathbf{w}\|_1 \leq \beta} J_S^{\text{LS}}(\mathbf{w})$$

- ▶ This is all to optimize the *expected future performance*.
- ▶ If we have infinite data, we don't need to worry about regularization.

General Definition of Overfitting

- ▶ **Population:** Input-output pairs (x, y) are distributed as \mathcal{D} .
- ▶ **Loss function:** A loss function $l(y, y') \in \mathbb{R}$ specifies the penalty on predicting y' when the correct answer is y .

General Definition of Overfitting

- ▶ **Population:** Input-output pairs (x, y) are distributed as \mathcal{D} .
- ▶ **Loss function:** A loss function $l(y, y') \in \mathbb{R}$ specifies the penalty on predicting y' when the correct answer is y .
- ▶ **Hypothesis class:** A class of functions \mathcal{H} is chosen to model the input-output relationship (e.g., all hyperplanes).

General Definition of Overfitting

- ▶ **Population:** Input-output pairs (x, y) are distributed as \mathcal{D} .
- ▶ **Loss function:** A loss function $l(y, y') \in \mathbb{R}$ specifies the penalty on predicting y' when the correct answer is y .
- ▶ **Hypothesis class:** A class of functions \mathcal{H} is chosen to model the input-output relationship (e.g., all hyperplanes).
- ▶ **Training data:** A fixed set of samples S is used to obtain your hypothesis

$$h_S = \arg \min_{h \in \mathcal{H}} \hat{\mathbf{E}}_S [l(y, h_S(x))] = \arg \min_{h \in \mathcal{H}} \frac{1}{|S|} \sum_{(x,y) \in S} l(y, h(x))$$

General Definition of Overfitting

- ▶ **Population:** Input-output pairs (x, y) are distributed as \mathcal{D} .
- ▶ **Loss function:** A loss function $l(y, y') \in \mathbb{R}$ specifies the penalty on predicting y' when the correct answer is y .
- ▶ **Hypothesis class:** A class of functions \mathcal{H} is chosen to model the input-output relationship (e.g., all hyperplanes).
- ▶ **Training data:** A fixed set of samples S is used to obtain your hypothesis

$$h_S = \arg \min_{h \in \mathcal{H}} \hat{\mathbf{E}}_S [l(y, h_S(x))] = \arg \min_{h \in \mathcal{H}} \frac{1}{|S|} \sum_{(x,y) \in S} l(y, h(x))$$

- ▶ We say h_S **overfits** S if there is $h \in \mathcal{H}$ such that

$$\begin{aligned} \hat{\mathbf{E}}_S [l(y, h_S(x))] &< \hat{\mathbf{E}}_S [l(y, h(x))] \\ \mathbf{E} [l(y, h_S(x))] &> \mathbf{E} [l(y, h(x))] \end{aligned}$$

Overview

Regularization

A Probabilistic Interpretation of Linear Regression

Optimization by Local Search

Linear Regression as MLE

- ▶ **Claim.** If $S = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$ is generated by a particular probabilistic model, then the least squares solution is also the maximum likelihood solution under this model:

$$\mathbf{w}_S^{\text{LS}} = \arg \max_{\mathbf{w} \in \mathbb{R}^d} \Pr(S|\mathbf{w})$$

Linear Regression as MLE

- ▶ **Claim.** If $S = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$ is generated by a particular probabilistic model, then the least squares solution is also the maximum likelihood solution under this model:

$$\mathbf{w}_S^{\text{LS}} = \arg \max_{\mathbf{w} \in \mathbb{R}^d} \Pr(S|\mathbf{w})$$

- ▶ Provides an alternative characterization of the method.
- ▶ This is a recurring theme: different approaches “converge” to the same thing.

The Probabilistic Model

- ▶ There is some input distribution $\mathbf{x} \sim \mathcal{D}$ (as before).
- ▶ The noise $\epsilon \sim \mathcal{N}(0, \sigma^2)$ is (centered) **Gaussian**.
- ▶ For any $\mathbf{w} \in \mathbb{R}^d$, the output value is set to $y = \mathbf{w} \cdot \mathbf{x} + \epsilon$.

The Probabilistic Model

- ▶ There is some input distribution $\mathbf{x} \sim \mathcal{D}$ (as before).
- ▶ The noise $\epsilon \sim \mathcal{N}(0, \sigma^2)$ is (centered) **Gaussian**.
- ▶ For any $\mathbf{w} \in \mathbb{R}^d$, the output value is set to $y = \mathbf{w} \cdot \mathbf{x} + \epsilon$.
- ▶ Thus

$$\begin{aligned}\Pr(\mathbf{x}, y | \mathbf{w}) &= \Pr(\mathbf{x}) \Pr(y | \mathbf{w}, \mathbf{x}) \\ &= \mathcal{D}(\mathbf{x}) \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-(y - \mathbf{w} \cdot \mathbf{x})^2}{2\sigma^2}\right) \right)\end{aligned}$$

MLE Coincides with Least Squares

Assuming each sample in $\mathcal{S} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$ is drawn iid,

$$\mathbf{w}_{\mathcal{S}}^{\text{MLE}} := \arg \max_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n \log \Pr(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w})$$

MLE Coincides with Least Squares

Assuming each sample in $S = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$ is drawn iid,

$$\begin{aligned}\mathbf{w}_S^{\text{MLE}} &:= \arg \max_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n \log \Pr(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) \\ &= \arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n \left(y^{(i)} - \mathbf{w} \cdot \mathbf{x}^{(i)} \right)^2\end{aligned}$$

MLE Coincides with Least Squares

Assuming each sample in $S = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$ is drawn iid,

$$\begin{aligned}\mathbf{w}_S^{\text{MLE}} &:= \arg \max_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n \log \Pr(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) \\ &= \arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n \left(y^{(i)} - \mathbf{w} \cdot \mathbf{x}^{(i)} \right)^2 = \mathbf{w}_S^{\text{LS}}\end{aligned}$$

Overview

Regularization

A Probabilistic Interpretation of Linear Regression

Optimization by Local Search

Need for a Universal Optimization Technique

- ▶ In the case of linear regression, the optimal model on a training dataset is given in closed-form by $\mathbf{w}_S^{\text{LS}} = X^+ \mathbf{y}$.

Need for a Universal Optimization Technique

- ▶ In the case of linear regression, the optimal model on a training dataset is given in closed-form by $w_S^{\text{LS}} = X^+y$.
- ▶ This **almost never happens** with a real world objective.

Need for a Universal Optimization Technique

- ▶ In the case of linear regression, the optimal model on a training dataset is given in closed-form by $w_S^{\text{LS}} = X^+y$.
- ▶ This **almost never happens** with a real world objective.

We need a **general**, **efficient** optimization technique that can be used for a wide class of models and objectives!

Local Search

Input: training objective* $J(\theta) \in \mathbb{R}$, number of iterations T

Output: parameter $\hat{\theta} \in \mathbb{R}^d$ such that $J(\hat{\theta})$ is small

1. Initialize θ^0 (e.g., randomly).
2. For $t = 0 \dots T - 1$,
 - 2.1 Obtain $\Delta^t \in \mathbb{R}^n$ such that $J(\theta^t + \Delta^t) \leq J(\theta^t)$.
 - 2.2 Choose some “step size” $\eta^t \in \mathbb{R}$.
 - 2.3 Set $\theta^{t+1} = \theta^t + \eta^t \Delta^t$.
3. Return θ^T .

* Assumed to be differentiable in this lecture.

Local Search

Input: training objective* $J(\theta) \in \mathbb{R}$, number of iterations T

Output: parameter $\hat{\theta} \in \mathbb{R}^d$ such that $J(\hat{\theta})$ is small

1. Initialize θ^0 (e.g., randomly).
2. For $t = 0 \dots T - 1$,
 - 2.1 Obtain $\Delta^t \in \mathbb{R}^n$ such that $J(\theta^t + \Delta^t) \leq J(\theta^t)$.
 - 2.2 Choose some “step size” $\eta^t \in \mathbb{R}$.
 - 2.3 Set $\theta^{t+1} = \theta^t + \eta^t \Delta^t$.
3. Return θ^T .

What is a good Δ^t ?

*Assumed to be differentiable in this lecture.

Gradient of the Objective at the Current Parameter

At $\theta^t \in \mathbb{R}^n$, the rate of increase (of the value of J) along a direction $u \in \mathbb{R}^n$ (i.e., $\|u\|_2 = 1$) is given by the **directional derivative**

$$\nabla_u J(\theta^t) := \lim_{\epsilon \rightarrow 0} \frac{J(\theta^t + \epsilon u) - J(\theta^t)}{\epsilon}$$

Gradient of the Objective at the Current Parameter

At $\theta^t \in \mathbb{R}^n$, the rate of increase (of the value of J) along a direction $u \in \mathbb{R}^n$ (i.e., $\|u\|_2 = 1$) is given by the **directional derivative**

$$\nabla_u J(\theta^t) := \lim_{\epsilon \rightarrow 0} \frac{J(\theta^t + \epsilon u) - J(\theta^t)}{\epsilon}$$

The **gradient** of J at θ^t is defined to be a vector $\nabla J(\theta^t)$ such that

$$\nabla_u J(\theta^t) = \nabla J(\theta^t) \cdot u \quad \forall u \in \mathbb{R}^n$$

Gradient of the Objective at the Current Parameter

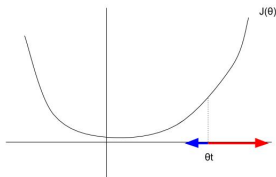
At $\theta^t \in \mathbb{R}^n$, the rate of increase (of the value of J) along a direction $u \in \mathbb{R}^n$ (i.e., $\|u\|_2 = 1$) is given by the **directional derivative**

$$\nabla_u J(\theta^t) := \lim_{\epsilon \rightarrow 0} \frac{J(\theta^t + \epsilon u) - J(\theta^t)}{\epsilon}$$

The **gradient** of J at θ^t is defined to be a vector $\nabla J(\theta^t)$ such that

$$\nabla_u J(\theta^t) = \nabla J(\theta^t) \cdot u \quad \forall u \in \mathbb{R}^n$$

Therefore, the **direction of the greatest rate of decrease** is given by $-\nabla J(\theta^t) / \|\nabla J(\theta^t)\|_2$.



Gradient Descent

Input: training objective $J(\theta) \in \mathbb{R}$, number of iterations T

Output: parameter $\hat{\theta} \in \mathbb{R}^n$ such that $J(\hat{\theta})$ is small

1. Initialize θ^0 (e.g., randomly).
2. For $t = 0 \dots T - 1$,

$$\theta^{t+1} = \theta^t - \eta^t \nabla J(\theta^t)$$

3. Return θ^T .

Gradient Descent

Input: training objective $J(\theta) \in \mathbb{R}$, number of iterations T

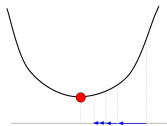
Output: parameter $\hat{\theta} \in \mathbb{R}^n$ such that $J(\hat{\theta})$ is small

1. Initialize θ^0 (e.g., randomly).
2. For $t = 0 \dots T - 1$,

$$\theta^{t+1} = \theta^t - \eta^t \nabla J(\theta^t)$$

3. Return θ^T .

When $J(\theta)$ is additionally *convex* (as in linear regression), gradient descent converges to an optimal solution (for appropriate step sizes).



Gradient Descent for Linear Regression

Input: training objective

$$J_S^{\text{LS}}(\mathbf{w}) := \frac{1}{2} \sum_{i=1}^n \left(y^{(i)} - \mathbf{w} \cdot \mathbf{x}^{(i)} \right)^2$$

number of iterations T

Output: parameter $\hat{\mathbf{w}} \in \mathbb{R}^n$ such that $J_S^{\text{LS}}(\hat{\mathbf{w}}) \approx J_S^{\text{LS}}(\mathbf{w}_S^{\text{LS}})$

1. Initialize \mathbf{w}^0 (e.g., randomly).
2. For $t = 0 \dots T - 1$,

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta^t \sum_{i=1}^n \mathbf{x}^{(i)} \cdot \left(y^{(i)} - \mathbf{w}^t \cdot \mathbf{x}^{(i)} \right)$$

3. Return \mathbf{w}^T .

Summary

- ▶ **Regularization** is an effort to prevent **overfitting** and optimize the true/population error.
- ▶ We can endow an alternative probabilistic interpretation of linear regression as MLE.
- ▶ **Gradient descent** is a local search algorithm that can be used to optimize *any* differentiable objective function.
 - ▶ A variant called “stochastic” gradient descent is the cornerstone of modern large-scale machine learning.