

# Multi-Agent Optimization and Learning

## Lecture 1: Foundations

Stefan Vlaski<sup>†</sup> and Ali H. Sayed<sup>\*</sup>

<sup>†</sup>Department of Electrical and Electronic Engineering, Imperial College London, UK  
<sup>\*</sup>Adaptive Systems Laboratory, École Polytechnique Fédérale de Lausanne, Switzerland

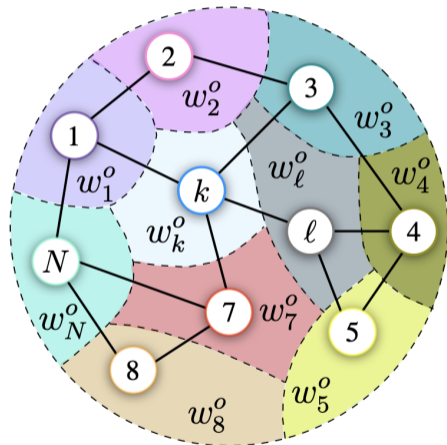
IEEE ICASSP 2024 Short Course

**Imperial College**  
**London**

**EPFL**

# Multi-Agent Systems with Dispersed Data and Capabilities

- Sensor networks
- Autonomous vehicles
- Social networks
- Power grids
- Mobile devices
- Drone swarms
- Many more



# Objectives of this Course

## Objective

To provide attendees with a fundamental understanding and the resulting intuition about distributed algorithms and the resulting performance trade-offs in intelligent multi-agent systems.

Hopefully, by the end of this course:

- **Those interested in developing distributed algorithms** will have the tools needed to rigorously develop and analyze new algorithms for distributed learning.
- **Those interested in applying distributed algorithms** will have an understanding of fundamental properties that allow for informed decisions on when and how to use and adapt distributed algorithms in their domain.

# Structure of this course

- **Lecture 1:** Foundations
  - ▶ Taxonomy, inference, learning and optimization, (stochastic) gradient descent.
- **Lecture 2:** Learning with a fusion center
  - ▶ Partial participation, local updates, heterogeneity and their impact on performance.
- **Lecture 3:** Graphs and their role decentralized processing
  - ▶ Graph spectral theory, averaging over graphs.
- **Lecture 4:** Algorithms for decentralized optimization and learning
  - ▶ Penalty-based (Distributed gradient descent and diffusion), primal-dual (EXTRA and Exact diffusion), and gradient-tracking (DIGing, NEXT and AugDGM)
- **Lecture 5:** Performance guarantees and trade-offs
  - ▶ Effect of heterogeneity and bias-correction, graph connectivity, gradient noise.
- **Lecture 6:** Advanced topics and open problems
  - ▶ Multi-task and meta-learning, compression, privacy.

# Outline of This Lecture

In this lecture we introduce some foundational concepts in a manner that will generalize and translate to distributed structures in the rest of the course.

- Components of a multi-agent system, and a taxonomy of distributed algorithms
- Inference, Learning and Optimization
- Deterministic and Stochastic Gradient Algorithms
- Centralized Learning

# Components of a Multi-Agent System

- **Agents:** Individual entity with basic capabilities to observe data, perform computations and communicate. Depending on the application, the term “agent” can represent a multitude of things, ranging from a sensor, to mobile devices such as phones, autonomous vehicles and drones, or humans.
- **Data:** Agents will be equipped with the ability to observe the world around them, and these observations will be captured in the form of data. These data may be available as a batch of fixed-size, or be observed on the fly. We will model data as random variables.
- **Objective:** Agents will aim to achieve some objective, which will be represented as an optimization problem.
- **Model:** In most applications, agents will maintain and refine a model, captured through parameters, which reconciles observed data and the objective. Models will be learned by adjusting the parameterization to optimize the agent’s objective.
- **Network:** We will represent communication links between agents through a network topology, which will place limits on the flow of information. We elaborate on some typical structures and taxonomy in the next slides.

## Non-cooperative learning

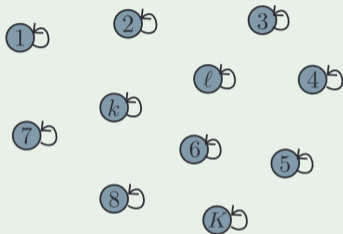
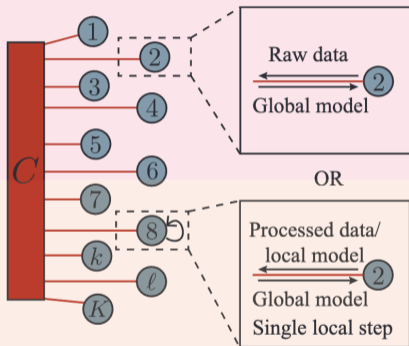


Figure: Non-cooperative approaches serve as a worst-case baseline.

## Fusion-center based learning

### Centralized



### Asynchronous or federated

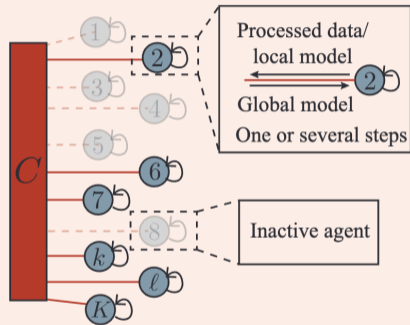


Figure: Centralized structures serve as a best-case baseline.



## Decentralized learning

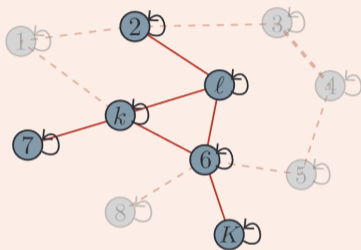


Figure: Decentralized structures avoid a fusion center.

# Maximum A Posteriori Estimation

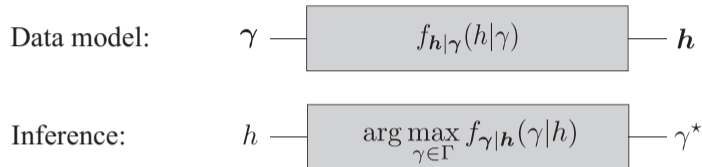
If we are provided with the conditional distribution  $f_{\gamma|\mathbf{h}}(\gamma|h)$  along with a single realization of the feature  $\mathbf{h} \in \mathbb{R}^{M_h}$ , it is quite natural to estimate  $\gamma \in \mathbb{R}^{M_\gamma}$  as the most likely outcome given  $\mathbf{h}$ . We can express this formally as:

$$\gamma^* \triangleq \arg \max_{\gamma \in \Gamma} f_{\gamma|\mathbf{h}}(\gamma|h) \quad (1)$$

The conditional distribution  $f_{\gamma|\mathbf{h}}(\gamma|h)$  is frequently referred to as the *posterior distribution* of  $\gamma$ , making  $\gamma^*$  the *maximum a posteriori (MAP) estimate* of  $\gamma$  given  $\mathbf{h}$ . We can also regard the MAP solution as a random variable, denoted in boldface notation  $\boldsymbol{\gamma}^*$ , when it is viewed as a function of the random observation  $\mathbf{h}$ :

$$\boldsymbol{\gamma}^* \triangleq \arg \max_{\gamma \in \Gamma} f_{\gamma|\mathbf{h}}(\gamma|\mathbf{h}) \quad (2)$$

# The Bayesian Learning Framework



**Figure:** A general model underlying most inference problems. The top row illustrates the generative model for the observations  $\mathbf{h}$ , while the bottom row illustrates the Bayesian inference formulation for recovering the label variable.

1

---

<sup>1</sup>We use **bold** font to describe a random variable, while regular font denotes its realization.

## Example: Linear Regression

Let us consider a feature  $\mathbf{h} \in \mathbb{R}^M$ , and a linear model parameterized by  $w^o \in \mathbb{R}^M$ , resulting in the label  $\gamma \in \mathbb{R}$ :

$$\gamma = \mathbf{h}^\top w^o + v \quad (3)$$

Here,  $v$  denotes zero-mean noise, which we assume to be normally distributed  $v \sim \mathcal{N}(0, \sigma_v^2)$  with variance  $\sigma_v^2 > 0$ . Suppose we would like to predict the most likely label  $\gamma$  for a given feature  $\mathbf{h}$ . We can first evaluate the conditional pdf:

$$f_{\gamma|\mathbf{h}}(\gamma|h) = \frac{1}{\sqrt{2\pi\sigma_v^2}} e^{-\frac{1}{2\sigma_v^2}(\gamma - \mathbf{h}^\top w^o)^2} \quad (4)$$

## Example: Linear Regression

We see directly that the posterior distribution  $f_{\gamma|h}(\gamma|h)$  is maximized whenever  $(\gamma - h^T w^o)^2 = 0$  and hence  $\gamma^* = h^T w^o$ . An alternative approach is to take the logarithm of the posterior and find:

$$\log f_{\gamma|h}(\gamma|h) = \log \left( \frac{1}{\sqrt{2\pi\sigma_v^2}} \right) - \frac{1}{2\sigma_v^2}(\gamma - h^T w^o)^2 \quad (5)$$

Since the logarithm of a positive variable is strictly increasing, it follows that the maximizing argument of  $f_{\gamma|h}(\gamma|h)$  is the same as the maximizing argument of  $\log f_{\gamma|h}(\gamma|h)$ , which is the choice  $\gamma$  that satisfies  $\frac{1}{2\sigma_v^2}(\gamma - h^T w^o)^2 = 0$ . It then follows that:

$$\gamma^* = \arg \max_{\gamma \in \Gamma} f_{\gamma|h}(\gamma|h) = h^T w^o \quad (6)$$

## Limitations

In its current form, however, the framework has two limitations.

- First, relation (1) and the following examples provided a rule to map a single observation  $h$  to an estimate  $\gamma^*$  of  $\gamma$ . In many inference and learning settings, we will be provided with multiple realizations  $\{h_n\}_{n=1}^N$  of the random variable  $\mathbf{h}$ , and seek to use these aggregate realizations in generating  $\gamma^*$ .
- Perhaps more importantly, the resulting estimation rules are driven in large part by the conditional probability densities (e.g.,  $f_{\gamma|\mathbf{h}}(\gamma|h)$  when estimating  $\gamma$  from  $\mathbf{h}$ ). This essentially requires full knowledge of the statistical model relating the label  $\gamma$  to its feature  $\mathbf{h}$ . For instance, in the case of the linear channel treated in the previous example, this translates into knowledge of the channel taps contained in  $w^o$ , along with the distribution of the noise term  $\mathbf{v}$ . In practice, we do not expect to have access to so much knowledge, and need to estimate relevant conditional distributions directly from data. We refer to the process of estimating models from data as “learning” and develop it from a Bayesian perspective.

## Multiple realizations

Dealing with multiple realizations is straightforward for i.i.d. data:

$$f_{\{\mathbf{h}_n\}_{n=1}^N|\gamma} \left( \{\mathbf{h}_n\}_{n=1}^N | \gamma \right) \stackrel{(a)}{=} \prod_{n=1}^N f_{\mathbf{h}_n|\gamma} (h_n|\gamma) \stackrel{(b)}{=} \prod_{n=1}^N f_{\mathbf{h}|\gamma} (h_n|\gamma) \quad (7)$$

where (a) holds by conditional independence and (b) holds by identical distribution of the random variables  $\mathbf{h}_n$ . We can then adjust the framework (1) to include multiple observations as follows:

$$\begin{aligned} \gamma^* &\triangleq \arg \max_{\gamma \in \Gamma} f_{\gamma|\{\mathbf{h}_n\}_{n=1}^N} \left( \gamma | \{\mathbf{h}_n\}_{n=1}^N \right) = \arg \max_{\gamma \in \Gamma} \frac{f_{\{\mathbf{h}_n\}_{n=1}^N|\gamma} \left( \{\mathbf{h}_n\}_{n=1}^N | \gamma \right) \times f_{\gamma}(\gamma)}{f_{\{\mathbf{h}_n\}_{n=1}^N} \left( \{\mathbf{h}_n\}_{n=1}^N \right)} \\ &= \arg \max_{\gamma \in \Gamma} f_{\{\mathbf{h}_n\}_{n=1}^N|\gamma} \left( \{\mathbf{h}_n\}_{n=1}^N | \gamma \right) \times f_{\gamma}(\gamma) \\ &= \arg \max_{\gamma \in \Gamma} \left( \prod_{n=1}^N f_{\mathbf{h}|\gamma} (h_n|\gamma) \right) \times f_{\gamma}(\gamma) \end{aligned} \quad (8)$$

## Multiple realizations

The fact that the logarithmic function is monotonically increasing ensures that:

$$\begin{aligned}\gamma^* &= \arg \max_{\gamma \in \Gamma} \log \left[ \left( \prod_{n=1}^N f_{\mathbf{h}|\gamma}(h_n|\gamma) \right) \times f_{\gamma}(\gamma) \right] \\ &= \arg \max_{\gamma \in \Gamma} \left\{ \sum_{n=1}^N \log f_{\mathbf{h}|\gamma}(h_n|\gamma) + \log f_{\gamma}(\gamma) \right\} \\ &= \arg \min_{\gamma \in \Gamma} \left\{ - \sum_{n=1}^N \log f_{\mathbf{h}|\gamma}(h_n|\gamma) - \log f_{\gamma}(\gamma) \right\}\end{aligned}\tag{9}$$

We can normalize by the sample size  $N$ , since the minimizing argument  $\gamma^*$  is unaffected by scaling:

$$\gamma^* = \arg \min_{\gamma \in \Gamma} \left\{ -\frac{1}{N} \sum_{n=1}^N \log f_{\mathbf{h}|\gamma}(h_n|\gamma) - \frac{1}{N} \log f_{\gamma}(\gamma) \right\}\tag{10}$$



## Example: Inverting a linear regression model

We consider a binary class variable  $\gamma \in \{+1, -1\}$ . Suppose that the observed feature  $\mathbf{h}$  is normally distributed with mean that depends on the class variable  $\gamma$ . Specifically, we model  $\mathbf{h} \sim \mathcal{N}(w^o, \sigma_v^2 I_M)$  when  $\gamma = +1$  and  $\mathbf{h} \sim \mathcal{N}(-w^o, \sigma_v^2 I_M)$  when  $\gamma = -1$ . This relation can be represented through the linear model:

$$\mathbf{h} = \gamma w^o + \mathbf{v} \quad (11)$$

with  $\mathbf{v} \sim \mathcal{N}(0, \sigma_v^2)$ . In this sense, we may view the present example as the inverse problem to the linear regression example we previously saw. The linear relationship translates into the conditional distribution:

$$f_{\mathbf{h}|\gamma}(\mathbf{h}|\gamma) = \frac{1}{\sqrt{(2\pi)^M \sigma_v^{2M}}} e^{-\frac{1}{2\sigma_v^2} \|\mathbf{h} - \gamma w^o\|^2} \quad (12)$$

## Example: Inverting a linear regression model

Applying (10), and assuming an equal prior  $\mathbb{P}\{\gamma = +1\} = \mathbb{P}\{\gamma = -1\} = \frac{1}{2}$ , we obtain:

$$\begin{aligned}\gamma^* &= \arg \min_{\gamma \in \pm 1} \frac{1}{N} \sum_{n=1}^N \|h_n - \gamma w^o\|^2 \\ &= \arg \min_{\gamma \in \pm 1} \left\{ \frac{1}{N} \sum_{n=1}^N \|h_n\|^2 - 2\gamma \left( \frac{1}{N} \sum_{n=1}^N h_n \right)^\top w^o + \|\gamma w^o\|^2 \right\} \\ &\stackrel{(a)}{=} \arg \min_{\gamma \in \pm 1} \left\{ -2\gamma \left( \frac{1}{N} \sum_{n=1}^N h_n \right)^\top w^o + \gamma^2 \|w^o\|^2 \right\} \\ &\stackrel{(b)}{=} \arg \min_{\gamma \in \pm 1} \left\{ -2\gamma \left( \frac{1}{N} \sum_{n=1}^N h_n \right)^\top w^o \right\} = \text{sign} \left\{ \left( \frac{1}{N} \sum_{n=1}^N h_n \right)^\top w^o \right\}\end{aligned}\quad (13)$$

In this sequence of reformulations, (a) follows from the fact that  $\frac{1}{N} \sum_{n=1}^N \|h_n\|^2$  is independent of  $\gamma$ , (b) uses the fact that  $\gamma^2 = 1$ .

## From Inference to Learning

So far, given a known model  $f_{\mathbf{h}|\gamma}$  and realizations of  $\mathbf{h}$ , we performed *inference* of the label  $\gamma$ . To formalize the *learning* process, we will assume that the conditional likelihood  $f_{\mathbf{h}|\gamma}(\cdot|\gamma)$  is parameterized by a learnable parameter  $w \in \mathbb{R}^{M_w}$  and write instead  $f_{\mathbf{h}|\gamma,w}(\cdot|\gamma,w)$ . Under this parameterization, learning the conditional distribution of  $\gamma$  given  $\mathbf{h}$  is equivalent to learning the parameter  $w$  that parameterizes  $f_{\mathbf{h}|\gamma,w}(\cdot|\gamma,w)$ .

To formulate a procedure for learning  $w$  from pairs  $\{\mathbf{h}, \gamma\}$  we mirror the argument so far. Suppose the model  $w$  is sampled once, yielding the realization  $w^o$ . The *training data*  $\{\mathbf{h}_n, \gamma_n\}_{n=1}^N$  is sampled  $N$  times, where each pair  $\{\mathbf{h}_n, \gamma_n\}$  is sampled from  $f_{\mathbf{h},\gamma|w}(h, \gamma|w^o)$ . If we suppose that feature-label pairs  $\{\mathbf{h}_n, \gamma_n\}$  are identically and independently distributed after conditioning on the parameterization  $w$ , we can factorize:

$$\begin{aligned} f_{\{\mathbf{h}_n, \gamma_n\}_{n=1}^N | w} \left( \{h_n, \gamma_n\}_{n=1}^N | w \right) &\stackrel{(a)}{=} \prod_{n=1}^N f_{\mathbf{h}_n, \gamma_n | w} (h_n, \gamma_n | w) \\ &\stackrel{(b)}{=} \prod_{n=1}^N f_{\mathbf{h}, \gamma | w} (h_n, \gamma_n | w) \end{aligned} \quad (14)$$

## From Inference to Learning

We can then define the MAP estimate of the weight vector  $w$  as:

$$\begin{aligned} w^* &\triangleq \arg \max_{w \in \mathbb{R}^{M_w}} f_{\mathbf{w}|\{\mathbf{h}_n, \gamma_n\}_{n=1}^N} \left( w | \{\mathbf{h}_n, \gamma_n\}_{n=1}^N \right) \\ &= \arg \max_{w \in \mathbb{R}^{M_w}} \frac{f_{\{\mathbf{h}_n, \gamma_n\}_{n=1}^N | \mathbf{w}} \left( \{\mathbf{h}_n, \gamma_n\}_{n=1}^N | w \right) \times f_{\mathbf{w}}(w)}{f_{\{\mathbf{h}_n, \gamma_n\}_{n=1}^N} \left( \{\mathbf{h}_n, \gamma_n\}_{n=1}^N \right)} \\ &= \arg \max_{w \in \mathbb{R}^{M_w}} f_{\{\mathbf{h}_n, \gamma_n\}_{n=1}^N | \mathbf{w}} \left( \{\mathbf{h}_n, \gamma_n\}_{n=1}^N | w \right) \times f_{\mathbf{w}}(w) \\ &= \arg \max_{w \in \mathbb{R}^{M_w}} \left( \prod_{n=1}^N f_{\mathbf{h}, \gamma | \mathbf{w}}(h_n, \gamma_n | w) \right) \times f_{\mathbf{w}}(w) \end{aligned} \quad (15)$$

Following the same argument that led to (10), we arrive at:

$$w^* = \arg \min_{w \in \mathbb{R}^M} \left\{ -\frac{1}{N} \sum_{n=1}^N \log f_{\mathbf{h}, \gamma | \mathbf{w}}(h_n, \gamma_n | w) - \frac{1}{N} \log f_{\mathbf{w}}(w) \right\} \quad (16)$$

## Distinction between the true model and the MAP estimate

Observe that we make a deliberate distinction between the *true* model  $w^o$ , which parameterizes the distribution  $f_{\mathbf{h}, \gamma | \mathbf{w}}(\mathbf{h}, \gamma | w^o)$  from which the samples  $\{\mathbf{h}_n, \gamma_n\}_{n=1}^N$  are generated, and the MAP *estimate*  $w^*$ , which maximizes the posterior distribution of  $\mathbf{w}$  after observing the samples  $\{\mathbf{h}_n, \gamma_n\}_{n=1}^N$ . In general, there will be a difference between the MAP model  $w^*$  and the true model  $w^o$ . The expectation is that as the size of the data set  $N$  grows, the MAP model  $w^*$  will become increasingly accurate and approach  $w^o$  as  $N \rightarrow \infty$ . Detailed proofs of generalization are beyond the scope of this course.

## Train and Test Data

Once we have found a  $w^*$  based on some training set  $\{h_n, \gamma_n\}_{n=1}^N$ , we can then apply the model to an unseen feature to perform *approximate* Bayesian inference:

$$\gamma^{\text{test}^*} \triangleq \arg \max_{\gamma \in \Gamma} f_{\gamma|h, w}(\gamma|h^{\text{test}}, w^*) \quad (17)$$

We emphasize that (17) is only an *approximate* MAP estimate for the true label  $\gamma^{\text{test}}$ , since it is generated using an estimate of the posterior distribution  $f_{\gamma|h, w}(\gamma|h^{\text{test}}, w^*)$  using the learned parameterization  $w^*$ . In general, and in particular for finite sample sizes  $N$  of the training data, the learned parameterization  $w^*$  will be different from the true parameterization  $w^o$  that actually generated the data. The difference between predictions made using the true model  $w^o$  and the learned model  $w^*$  is known as a *generalization error*.

## Compact notation for feature-label pairs

As is evident from (16), during learning, we are provided with feature-label pairs  $\{\mathbf{h}, \gamma\}$  or their realizations  $\{h_n, \gamma_n\}_{n=1}^N$ . To simplify the notation, we will collect features  $\mathbf{h} \in \mathbb{R}^{M_h}$  and labels  $\gamma \in \mathbb{R}^{M_\gamma}$  into a single augmented data vector  $\mathbf{x} \in \mathbb{R}^{M_h + M_\gamma}$ , such that:

$$\mathbf{x} \triangleq \text{col} \{\mathbf{h}, \gamma\} = \begin{pmatrix} \mathbf{h} \\ \gamma \end{pmatrix} \quad (18)$$

Similarly, we will collect realizations into  $x_n \triangleq \text{col} \{h_n, \gamma_n\}$ . In this manner, we can write more compactly:

$$f_{\mathbf{h}, \gamma | \mathbf{w}}(h, \gamma | w) = f_{\mathbf{x} | \mathbf{w}}(x | w) \quad (19)$$

The MAP learning problem (16) then becomes:

$$w^* = \arg \min_{w \in \mathbb{R}^M} \left\{ -\frac{1}{N} \sum_{n=1}^N \log f_{\mathbf{x} | \mathbf{w}}(x_n | w) - \frac{1}{N} \log f_{\mathbf{w}}(w) \right\} \quad (20)$$

## Ridge regression for linear models

Let us consider again the linear model from previous examples.

$$\gamma = \mathbf{h}^\top \mathbf{w}^o + v \quad (21)$$

We can interpret  $w^o$  as the realization of a random variable  $\mathbf{w}$ , and assume that the parameter  $\mathbf{w}$  is independent of all other random variables. Suppose we impose a normal prior of the form  $\mathbf{w} \sim \mathcal{N}(0, \sigma_w^2 I_M)$ . We continue to assume  $v \sim \mathcal{N}(0, \sigma_v^2)$  and  $\mathbf{h} \sim \mathcal{N}(0, \sigma_h^2 I_M)$ . It then follows that:

$$f_{\mathbf{w}}(\mathbf{w}) = \frac{1}{\sqrt{(2\pi)^M \sigma_w^{2M}}} e^{-\frac{1}{2\sigma_w^2} \|\mathbf{w}\|^2} \quad (22)$$

$$\begin{aligned} f_{\mathbf{h}, \gamma | \mathbf{w}}(h, \gamma | \mathbf{w}) &= f_{\gamma | \mathbf{h}, \mathbf{w}}(\gamma | h, \mathbf{w}) \times f_{\mathbf{h} | \mathbf{w}}(h | \mathbf{w}) \\ &\stackrel{(a)}{=} f_{\gamma | \mathbf{h}, \mathbf{w}}(\gamma | h, \mathbf{w}) \times f_{\mathbf{h}}(h) \\ &= \frac{1}{\sqrt{2\pi\sigma_v^2}} e^{-\frac{1}{2\sigma_v^2} (\gamma - h^\top \mathbf{w})^2} \times \frac{1}{\sqrt{(2\pi)^M \sigma_h^{2M}}} e^{-\frac{1}{2\sigma_h^2} \|h\|^2} \end{aligned} \quad (23)$$



## Ridge regression for linear models

Substituting these relations into (16), and removing constant terms independent of  $w$ , we obtain:

$$w^* = \arg \min_{w \in \mathbb{R}^M} \left\{ \frac{1}{2N} \sum_{n=1}^N \left( \gamma_n - h_n^\top w \right)^2 + \frac{1}{2N} \frac{\sigma_v^2}{\sigma_w^2} \|w\|^2 \right\} \quad (24)$$

which corresponds to a least-squares regression problem, where the fit of the parameter  $w$  is evaluated by measuring the squared error of the prediction  $h_n^\top w$  to the observation  $\gamma_n$ . In practice, the variances  $\sigma_v^2$  and  $\sigma_w^2$  are unknown, and it is common to instead introduce a *hyperparameter*  $\rho$  resulting in:

$$w^* = \arg \min_{w \in \mathbb{R}^M} \left\{ \frac{1}{N} \sum_{n=1}^N \left( \gamma_n - h_n^\top w \right)^2 + \rho \|w\|^2 \right\} \quad (25)$$

The hyperparameter  $\rho$  is either fixed a priori based on an estimate of the ratio  $\frac{\sigma_v^2}{\sigma_w^2}$ , or tuned on a subset of the training data using cross-validation.

## Asymptotic Behavior

It is instructive to examine the asymptotic behavior of the optimal estimate (16) as the sample size  $N$  grows. We have for all  $w$  such that  $f_w(w) > 0$ :

$$\begin{aligned} & \lim_{N \rightarrow \infty} \left\{ -\frac{1}{N} \sum_{n=1}^N \log f_{\mathbf{x} | \mathbf{w}}(x_n | w) - \frac{1}{N} \log f_w(w) \right\} \\ &= \lim_{N \rightarrow \infty} \left\{ -\frac{1}{N} \sum_{n=1}^N \log f_{\mathbf{x} | \mathbf{w}}(x_n | w) \right\} - \lim_{N \rightarrow \infty} \left\{ \frac{1}{N} \log f_w(w) \right\} \\ &\stackrel{(a)}{=} \lim_{N \rightarrow \infty} \left\{ -\frac{1}{N} \sum_{n=1}^N \log f_{\mathbf{x} | \mathbf{w}}(x_n | w) \right\} \\ &\stackrel{(b)}{=} -\mathbb{E}_{\mathbf{x} \sim f_{\mathbf{x} | \mathbf{w}}(\mathbf{x} | w)} \left\{ \log f_{\mathbf{x} | \mathbf{w}}(\mathbf{x} | w) \right\} \end{aligned} \tag{26}$$

Here, (a) follows since  $f_w(w) > 0 \implies |\log f_w(w)| < \infty$  and hence:

$$\lim_{N \rightarrow \infty} \left\{ \frac{1}{N} \log f_w(w) \right\} = 0 \tag{27}$$

## Example: Mean-Square-Error (MSE) estimation

Considering the linear model from previous examples, we find from (23):

$$\begin{aligned} & -\log f_{\mathbf{h}, \gamma | \mathbf{w}}(\mathbf{h}, \gamma | \mathbf{w}) \\ \stackrel{(23)}{=} & -\log \left( \frac{1}{\sqrt{2\pi\sigma_v^2}} e^{-\frac{1}{2\sigma_v^2}(\gamma - \mathbf{h}^\top \mathbf{w})^2} \times \frac{1}{\sqrt{(2\pi)^M \sigma_h^{2M}}} e^{-\frac{1}{2\sigma_h^2} \|\mathbf{h}\|^2} \right) \\ = & -\log \left( e^{-\frac{1}{2\sigma_v^2}(\gamma - \mathbf{h}^\top \mathbf{w})^2} \right) - \underbrace{\log \left( \frac{1}{\sqrt{2\pi\sigma_v^2}} \times \frac{1}{\sqrt{(2\pi)^M \sigma_h^{2M}}} e^{-\frac{1}{2\sigma_h^2} \|\mathbf{h}\|^2} \right)}_{\triangleq C} \\ = & \frac{1}{2\sigma_v^2} (\gamma - \mathbf{h}^\top \mathbf{w})^2 + C \end{aligned} \tag{28}$$

Since  $\frac{1}{\sigma_v^2}$  and  $C$  are independent of  $w$ , we conclude from (26) after shifting and scaling:

$$w^o = \arg \min_w \frac{1}{2} \mathbb{E} \left( \gamma - \mathbf{h}^\top \mathbf{w} \right)^2 \tag{29}$$

## Example: Mean-Square-Error (MSE) estimation

Since the risk function in this case is quadratic in  $w$ , we can obtain a closed-form expression for the minimizer  $w^o$ . Indeed, after differentiating, we find:

$$\nabla \left( \frac{1}{2} \mathbb{E} \left( \gamma - \mathbf{h}^T w \right)^2 \right) = - \underbrace{\mathbb{E} \gamma \mathbf{h}}_{\triangleq r_{\gamma h}} + \underbrace{\mathbb{E} \mathbf{h} \mathbf{h}^T}_{\triangleq R_h} w = -r_{\gamma h} + R_h w \quad (30)$$

At  $w^o$ , this derivative must be equal to zero, and hence after rearranging:

$$w^o = R_h^{-1} r_{\gamma h} \quad (31)$$

We can verify that the minimizing argument  $w^o$  of the MSE risk (29) indeed corresponds to the *true* model parameters in (3). To this end, note from

$$\gamma = \mathbf{h}^T w^o + \mathbf{v} \quad (32)$$

that if we multiply by  $\mathbf{h}$  from the left and take expectations, we get:

$$r_{\gamma h} = R_h w^o \quad (33)$$

# Empirical and expected risk minimization for single agent learning

- We have hence motivated the following optimization problems as rigorous learning frameworks.
- General empirical risk minimization problems take the form:

$$w^{\star} = \arg \min_w \frac{1}{N} \sum_{n=1}^N Q(w; x_n) + R(w) \quad (34)$$

Here,  $Q(w; x_n)$  denotes the *loss* of the parameterization  $w$  on the sample  $x_n$ . The term  $R(w)$  corresponds to a regularization term, which is independent of the data  $x_n$ , but may depend on the sample size  $N$ .

- (Expected) risk minimization problems take the form:

$$w^o = \arg \min_w \mathbb{E}_{\mathbf{x}} Q(w; \mathbf{x}) \quad (35)$$

## Extension to Multi-Agent Systems

To this end, we consider a collection of  $K$  agents. Each agent has sensing capabilities, and collects a local dataset  $\{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N$ . For simplicity, we assume that the sample size  $N$  is common to all agents. Each pair of samples is sampled independently and identically from the distribution:

$$f_{\mathbf{h}_{k,n}, \gamma_{k,n} | \mathbf{w}_k} (h, \gamma | w^o) = f_{\mathbf{h}, \gamma | \mathbf{w}} (h, \gamma | w^o) \quad (36)$$

Now assume first that the different agents are disconnected, with no way to exchange information among them. In that case, it is reasonable to formulate independent local learning problems:

$$w_k^* \triangleq \arg \max_{w \in \mathbb{R}^M} f_{\mathbf{w} | \{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N} (w | \{h_{k,n}, \gamma_{k,n}\}_{n=1}^N) \quad (37)$$

## Extension to Multi-Agent Systems

Following the argument same argument as before, we arrive at the equivalent formulation:

$$\begin{aligned} w_k^* &= \arg \min_{w \in \mathbb{R}^M} \left\{ -\frac{1}{N} \sum_{n=1}^N \log f_{\mathbf{h}, \gamma | w} (h_{k,n}, \gamma_{k,n} | w) - \frac{1}{N} \log f_w(w) \right\} \\ &= \arg \min_{w \in \mathbb{R}^M} J_k(w) \end{aligned} \quad (38)$$

where we defined:

$$J_k(w) \triangleq -\frac{1}{N} \sum_{n=1}^N \log f_{\mathbf{h}, \gamma | w} (h_{k,n}, \gamma_{k,n} | w) - \frac{1}{N} \log f_w(w) \quad (39)$$

Through the agent subscript  $k$  we emphasize here that because the local samples  $\{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N$  are distinct, this implies that  $J_k(\cdot)$  as well as the resulting optimal models  $w_k^*$  are distinct as well. While we can claim local optimality of each model  $w_k^*$  in the sense of (38), this model is making use only of the data available to agent  $k$ .

## Extension to Multi-Agent Systems

As an alternative to (38), we can instead consider the *global* learning problem:

$$w^* \triangleq \arg \max_{w \in \mathbb{R}^M} f_{\mathbf{w} | \left\{ \left\{ \mathbf{h}_{k,n}, \gamma_{k,n} \right\}_{n=1}^N \right\}_{k=1}^K} \left( w | \left\{ \left\{ h_{k,n}, \gamma_{k,n} \right\}_{n=1}^N \right\}_{k=1}^K \right) \quad (40)$$

As long as the samples  $\left\{ \mathbf{h}_{k,n}, \gamma_{k,n} \right\}_{n=1}^N$  are also mutually independent between agents after conditioning on  $w$ , we can factorize the conditional likelihoods:

$$\begin{aligned} w^* &\triangleq \arg \max_{w \in \mathbb{R}^M} f_{\mathbf{w} | \left\{ \left\{ \mathbf{h}_{k,n}, \gamma_{k,n} \right\}_{n=1}^N \right\}_{k=1}^K} \left( w | \left\{ \left\{ h_{k,n}, \gamma_{k,n} \right\}_{n=1}^N \right\}_{k=1}^K \right) \\ &= \arg \max_{w \in \mathbb{R}^M} f_{\left\{ \left\{ \mathbf{h}_{k,n}, \gamma_{k,n} \right\}_{n=1}^N \right\}_{k=1}^K | \mathbf{w}} \left( \left\{ \left\{ h_{k,n}, \gamma_{k,n} \right\}_{n=1}^N \right\}_{k=1}^K | w \right) \times f_{\mathbf{w}}(w) \\ &= \arg \max_{w \in \mathbb{R}^M} \left( \prod_{k=1}^K f_{\left\{ \mathbf{h}_{k,n}, \gamma_{k,n} \right\}_{n=1}^N | \mathbf{w}} \left( \left\{ h_{k,n}, \gamma_{k,n} \right\}_{n=1}^N | w \right) \right) \times f_{\mathbf{w}}(w) \end{aligned} \quad (41)$$



# Aggregate Optimization Problem

After taking logarithms, and normalizing, we find:

$$\begin{aligned} w^* &= \arg \min_w \left\{ -\frac{1}{KN} \sum_{k=1}^K \sum_{n=1}^N \log f_{h,\gamma|w} (h_{k,n}, \gamma_{k,n}|w) - \frac{1}{KN} \log f_w(w) \right\} \\ &= \arg \min_w \frac{1}{K} \sum_{k=1}^K J_k(w) \\ &= J(w) \end{aligned} \tag{42}$$

where we defined:

$$J(w) = \frac{1}{K} \sum_{k=1}^K J_k(w) \tag{43}$$

This problem is known as an *aggregate optimization problem* or *consensus optimization problem* for the collection of agents  $1, \dots, K$ .

## Example: Least-Squares Fusion

We specialize the discussion again to the least-squares problem. For the local models we would have:

$$w_k^* = \arg \min_{w \in \mathbb{R}^M} \left\{ \frac{1}{2N} \sum_{n=1}^N \left( \gamma_{k,n} - h_{k,n}^\top w \right)^2 + \frac{\rho}{2} \|w\|^2 \right\} \quad (44)$$

where we defined  $\rho = \frac{1}{N} \frac{\sigma_v^2}{\sigma_w^2}$  for notational convenience. Since (44) is quadratic over  $w$ , we can determine a closed-form expression for  $w_k^*$  by differentiating:

$$-\frac{1}{N} \sum_{n=1}^N \gamma_{k,n} h_{k,n} + \left( \frac{1}{N} \sum_{n=1}^N h_{k,n} h_{k,n}^\top \right) w_k^* + \rho w_k^* = 0 \quad (45)$$

## Example: Least-Squares Fusion

After grouping terms:

$$\begin{aligned} w_k^* &= \left( \frac{1}{N} \sum_{n=1}^N h_{k,n} h_{k,n}^\top + \rho I \right)^{-1} \left( \frac{1}{N} \sum_{n=1}^N \gamma_{k,n} h_{k,n} \right) \\ &= (H_k + \rho I)^{-1} d_k \end{aligned} \quad (46)$$

We defined for compactness:

$$H_k \triangleq \frac{1}{N} \sum_{n=1}^N h_{k,n} h_{k,n}^\top \quad (47)$$

$$d_k \triangleq \frac{1}{N} \sum_{n=1}^N \gamma_{k,n} h_{k,n} \quad (48)$$

## Example: Least-Squares Fusion

For the aggregate problem (43) we can find analogously:

$$\begin{aligned} w^* &= \left( \frac{1}{KN} \sum_{k=1}^K \sum_{n=1}^N h_{k,n} h_{k,n}^\top + \frac{\rho}{K} I \right)^{-1} \left( \frac{1}{KN} \sum_{k=1}^K \sum_{n=1}^N \gamma_{k,n} h_{k,n} \right) \\ &= \left( H + \frac{\rho}{K} I \right)^{-1} d \end{aligned} \quad (49)$$

where we defined:

$$H = \frac{1}{KN} \sum_{k=1}^K \sum_{n=1}^N h_{k,n} h_{k,n}^\top = \frac{1}{K} \sum_{k=1}^K H_k \quad (50)$$

$$d = \frac{1}{KN} \sum_{k=1}^K \sum_{n=1}^N \gamma_{k,n} h_{k,n} = \frac{1}{K} \sum_{k=1}^K d_k \quad (51)$$

## Exchanging Processed Quantities

Returning to the quadratic problem, note that:

$$\left(H + \frac{\rho}{K}I\right)w^* = d = \frac{1}{K} \sum_{k=1}^K d_k \stackrel{(46)}{=} \frac{1}{K} \sum_{k=1}^K (H_k + \rho I)w_k^* \quad (52)$$

It then follows that  $w^*$  can be evaluated equivalently as:

$$w^* = \left(\frac{1}{K} \sum_{k=1}^K H_k + \frac{\rho}{K}I\right)^{-1} \left(\frac{1}{K} \sum_{k=1}^K (H_k + \rho I)w_k^*\right) \quad (53)$$

This means that rather than communicate raw data  $h_{k,n}$  at the fusion center, and subsequently compute  $w^*$  via (49), we can instead locally compute  $H_k$  and  $w_k^*$ , and transmit only these processed quantities to the fusion center, and then apply (53).

- The first approach requires the communication of  $K(M+1)N$  scalar quantities.
- The second approach requires  $K$  times the exchange of a matrix of size  $M^2$  and a vector of size  $M$ , hence a cost of  $K(M+1)M$ .

## Recap and Outlook

- We used Bayesian techniques to motivate empirical and expected risk minimization problems of the form:

$$w^* = \arg \min_w \sum_{n=1}^N Q(w; x_n) + R(w) \triangleq \arg \min_w J(w) \quad (54)$$

$$w^o = \arg \min_w \mathbb{E}_{\mathbf{x}} Q(w; \mathbf{x}) \triangleq \arg \min_w J(w) \quad (55)$$

- We also saw how same line of reasoning motivates **aggregate** or **consensus** optimization problems of the form:

$$J(w) = \frac{1}{K} \sum_{k=1}^K J_k(w) \quad (56)$$

- In the case of linear models with Gaussian noise, this led to quadratic problems where we could do all calculations in closed-form.
- For more general objective functions  $J(w)$ , we will have to appeal to iterative methods.

# Deterministic Optimization

We are interested in optimizing a general objective:

$$w^* \triangleq \arg \min_w J(w) \quad (57)$$

- We need to assume some regularity conditions. For the purposes of this course, we will focus on  $\nu$ -strongly convex objectives with  $\delta$ -Lipschitz gradients.
- It is possible to relax both the convexity and smoothness conditions on the risks, but many of the theorems and conclusions will continue to hold. So to avoid unnecessary complications, and focus on key ideas, we will work with convex and smooth objectives.

## $\nu$ -Strong Convexity

We will say that the risk is  $\nu$ -strongly convex if:

$$J(w_2) \geq J(w_1) + \nabla J(w)^\top (w_2 - w_1) + \frac{\nu}{2} \|w_2 - w_1\|^2 \quad (58)$$

for some positive parameter  $\nu \geq 0$ . Note that strong convexity is a stronger condition than convexity, and in particular we recover convexity when  $\nu = 0$ . Now suppose we have found, by whatever means, a point satisfying the first-order condition  $\nabla J(w^\star) = 0$ . Then, it follows from (58) that for any  $w \in \text{dom } J$ :

$$\begin{aligned} J(w) &\geq J(w^\star) + \nabla J(w^\star)^\top (w - w^\circ) + \frac{\nu}{2} \|w - w^\circ\|^2 \\ &= J(w^\circ) + \frac{\nu}{2} \|w - w^\circ\|^2 \end{aligned} \quad (59)$$

Hence,  $w^\circ$  is the unique minimizing argument of  $J(w)$ .



## $\delta$ -Lipschitz Gradients

We will say that  $J(w)$  is smooth if it satisfies:

$$J(w_2) \leq J(w_1) + \nabla J(w_1)^\top (w_2 - w_1) + \frac{\delta}{2} \|w_2 - w_1\|^2 \quad (60)$$

Condition (60) is analogous to (58), except that now we impose an *upper* bound on the growth above the tangent hyperplane, where  $\delta$  controls the rate of growth. This implies the statement:

$$\|\nabla J(w_2) - \nabla J(w_1)\| \leq \delta \|w_2 - w_1\| \quad (61)$$

In other words, the smoothness condition (60) translates into a Lipschitz condition on the gradients of  $J(w)$ .

# Gradient Descent

We can then pursue

$$w^* \triangleq \arg \min_w J(w) \quad (62)$$

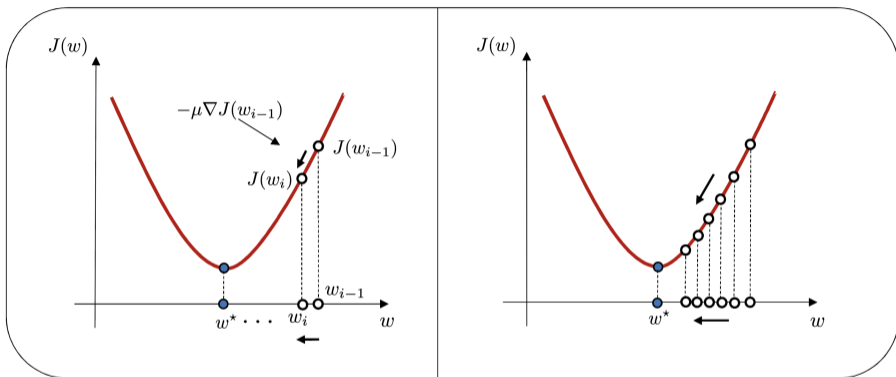
by iterating from some initial point  $w_0 \in \mathbb{R}^M$ :

$$w_i = w_{i-1} - \mu \nabla J(w_{i-1}) \quad (63)$$

which is known as the gradient descent algorithm. The name comes from the fact that:

$$\begin{aligned} J(w_i) &\stackrel{(60)}{\leq} J(w_{i-1}) + \nabla J(w_{i-1})^\top (w_i - w_{i-1}) + \frac{\delta}{2} \|w_i - w_{i-1}\|^2 \\ &\stackrel{(63)}{=} J(w_{i-1}) + \nabla J(w_{i-1})^\top (-\mu \nabla J(w_{i-1})) + \frac{\delta}{2} \|\mu \nabla J(w_{i-1})\|^2 \\ &= J(w_{i-1}) - \mu \left(1 - \mu \frac{\delta}{2}\right) \|\nabla J(w_{i-1})\|^2 \\ &\stackrel{(a)}{\leq} J(w_{i-1}) \end{aligned} \quad (64)$$

# Visualization of Gradient Descent



**Figure:** The panel on the left shows the mechanics of one update step where  $w_{i-1}$  is updated in the direction of the minimizer  $w^*$ . The panel on the right shows the result of several successive steps with the iterates approaching  $w^*$ .

## Example: Logistic Regression

Quadratic loss functions are not appropriate when the label  $\gamma$  in an inference problem takes discrete values  $\gamma = \pm 1$ . In this case, a more appropriate choice is the regularized logistic risk function defined by

$$J(w) = \frac{\rho}{2} \|w\|^2 + \frac{1}{N} \sum_{n=1}^N \ln \left( 1 + e^{-\gamma_n h_n^\top w} \right) \quad (65)$$

for which

$$\nabla J(w) = \rho w - \frac{1}{N} \sum_{n=1}^N \frac{\gamma_n h_n}{1 + e^{\gamma_n h_n^\top w}} \quad (66)$$

## Example: Logistic Regression

The gradient-descent recursion (63) in this case leads to

$$w_i = (1 - \mu\rho) w_{i-1} + \mu \left( \frac{1}{N} \sum_{n=1}^N \frac{\gamma_n h_n}{1 + e^{\gamma_n h_n^\top w_{i-1}}} \right), \quad i \geq 0 \quad (67)$$

One can interpret the logistic risk function as arising from a Bayes' optimal estimate for  $w$  under the statistical model:

$$f_{\gamma|h,w}(\gamma|h, w) = \frac{1}{1 + e^{-\gamma h^\top w}} \quad (68)$$

with a Gaussian prior on  $w$ , though the details of this derivation are not relevant to this course. Once we have found  $w^*$ , we can then perform inference on unlabelled data  $h$  by computing  $\text{sign}(h^\top w^*)$  as illustrated on the next slide. Features lying on one side of the hyperplane whose normal direction is  $w^*$  belong to class +1 while the other features belong to class -1.

## Example: Logistic Regression

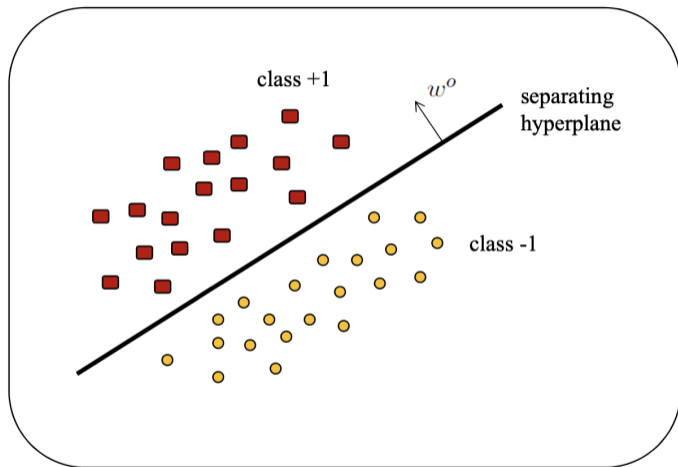
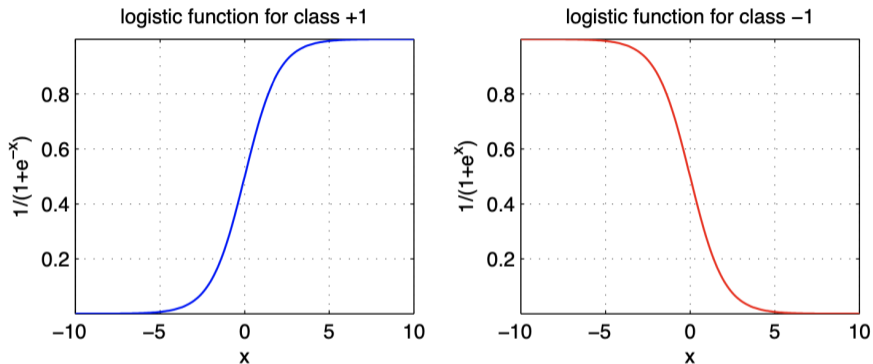


Figure: A logistic regression problem, with separating hyperplane defined by the parameter  $w^*$ .

## Example: Logistic Regression



**Figure:** The left figure shows the logistic function (68) when  $\gamma = +1$  as a function of  $x = h^T w$ , the right figure shows the same for  $\gamma = -1$ . This model imposes that it is likely that  $\gamma$  and  $h^T w$  have the same sign, with the likelihood increasing with the absolute value of  $h^T w$ .

# Convergence Analysis

## Convergence Of Gradient-Descent

Consider the gradient-descent recursion (63) where  $J(w)$  is a  $\nu$ -strongly convex function with  $\delta$ -Lipschitz gradients. Introduce the error vector  $\tilde{w}_i = w^* - w_i$ , which measures the difference between the  $i$ -th iterate and the global minimizer of  $J(w)$ . If the step-size  $\mu$  satisfies:

$$0 < \mu < 2\nu/\delta^2 \quad (69)$$

then  $w_i$  converges exponentially fast to  $w^*$  in the sense that

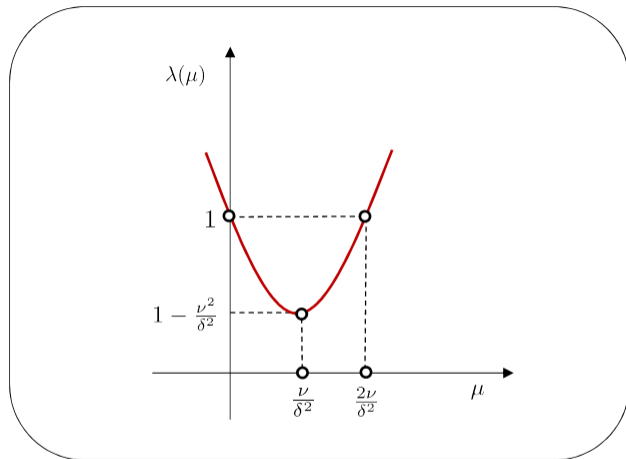
$$\|\tilde{w}_i\|^2 \leq \lambda \|\tilde{w}_{i-1}\|^2, \quad i \geq 0 \quad (70)$$

where  $\lambda = 1 - 2\mu\nu + \mu^2\delta^2 \in [0, 1)$ . It also holds that the risk value converges exponentially fast as follows

$$J(w_i) - J(w^*) \leq \delta\lambda^i \|\tilde{w}_0\|^2 = O(\lambda^i) \quad (71)$$

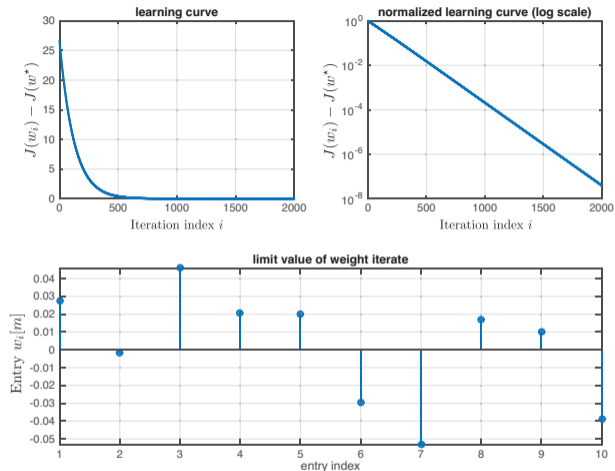


# Rate of Convergence



**Figure:** Plot of the function  $\lambda(\mu) = 1 - 2\nu\mu + \mu^2\delta^2$ . It shows that the function  $\lambda(\mu)$  assumes values below one in the range  $0 < \mu < 2\nu/\delta^2$ .

# Transient and steady-state error



**Figure:** Learning curves  $J(w_i)$  relative to the minimum risk value  $J(w^*)$  in linear scale (on the left) and in normalized logarithmic scale (on the right). (*Bottom*) Limiting value of  $w_i$ , which tends to  $w^*$ .

# Recap and Outlook

- We have seen how the gradient algorithm can iteratively minimize general objective functions.
- However, the recursions relied on exact gradients. In practice, this can be unavailable or costly, particularly in many applications of multi-agent systems.
- We will now see how we can tackle these issues with stochastic gradient algorithms.

# Motivation for Stochastic Gradients

We will continue to study optimization problems of the form:

$$\min_{w \in \mathbb{R}^M} J(w) \quad (72)$$

In the previous section, we examined the performance of the gradient-descent algorithm, which takes the form:

$$w_i = w_{i-1} - \mu \nabla J(w_{i-1}) \quad (73)$$

In practice, many times computing  $\nabla J(w_{i-1})$  is costly or infeasible.

## Example: Expected Risk Minimization

Suppose the risk  $J(w)$  is defined as the expected value of some loss function  $Q(w; \mathbf{x})$  over the distribution of the data  $\mathbf{x}$ :

$$J(w) \triangleq \mathbb{E}_{\mathbf{x}} Q(w; \mathbf{x}) \quad (74)$$

It is clear that the evaluation of  $J(w)$  requires knowledge of the distribution of  $\mathbf{x}$ , as well as the evaluation of its gradient vector since:

$$\nabla J(w) = \nabla (\mathbb{E}_{\mathbf{x}} Q(w; \mathbf{x})) \quad (75)$$

Suppose at time  $i$ , we observe a single sample  $\mathbf{x}_i$  from the distribution of  $\mathbf{x}$ . Then, we can construct the following instantaneous approximation for the gradient vector:

$$\widehat{\nabla J}(w; \mathbf{x}_i) \triangleq \nabla Q(w; \mathbf{x}_i) \quad (76)$$

and use it in (73):

$$\mathbf{w}_i = \mathbf{w}_{i-1} - \mu \widehat{\nabla J}(\mathbf{w}_{i-1}; \mathbf{x}_i) = \mathbf{w}_{i-1} - \mu \nabla Q(\mathbf{w}_{i-1}; \mathbf{x}_i) \quad (77)$$

Observe that the stochastic gradient approximation  $\nabla Q(\mathbf{w}_{i-1}; \mathbf{x}_i)$  is a function of the random variable  $\mathbf{x}_i$ , and is hence random itself.

## Example: Empirical Risk Minimization

A second class of problems that appear in learning contexts are empirical risk minimization problems of the form:

$$J(w) \triangleq \frac{1}{N} \sum_{n=1}^N Q(w; x_n) \quad (78)$$

where the set of points  $\{x_n\}_{n=1}^N$  corresponds to a batch of  $N$  independent realizations of the random variable  $\mathbf{x}$ . In contrast to the previous example, where evaluation of the gradient  $\nabla J(w)$  is infeasible, for empirical risk minimization problems, it is in principle possible to iterate:

$$w_i = w_{i-1} - \mu \nabla J(w_{i-1}) = w_{i-1} - \frac{\mu}{N} \sum_{n=1}^N \nabla Q(w_{i-1}; x_n) \quad (79)$$

The drawback now is that every iteration requires the evaluation of  $N$  gradients of the risk  $\nabla Q(w_{i-1}; x_n)$ . This results in a costly procedure, especially in large-scale settings, where the sample size  $N$  is large.

## Example: Empirical Risk Minimization

We may motivate a stochastic approximation for the true gradient appearing in (79) as follows. At any time  $i$ , we sample from the set of data points  $\{x_n\}_{n=1}^N$  uniformly at random. We may model this procedure by defining a random index  $n_i$ , with uniform distribution:

$$n_i = \begin{cases} 1, & \text{with probability } \frac{1}{N}, \\ 2, & \text{with probability } \frac{1}{N}, \\ \vdots & \\ N, & \text{with probability } \frac{1}{N}. \end{cases} \quad (80)$$

Then  $x_{n_i}$  denotes the sample picked randomly from  $\{x_n\}_{n=1}^N$  at time  $i$ . We use this sample to construct a stochastic approximation for the gradient as

$$\widehat{\nabla} J(w; x_{n_i}) \triangleq \nabla Q(w; x_{n_i}) \quad (81)$$

This construction results in a stochastic variant of recursion (79):

$$w_i = w_{i-1} - \mu \nabla Q(w_{i-1}, x_{n_i}) \quad (82)$$

## Example: Empirical Risk Minimization

Note further that if we let  $\mathbf{x}_i^{\text{emp}} \triangleq \mathbf{x}_{n_i}$ , it follows that:

$$J(w) \triangleq \frac{1}{N} \sum_{n=1}^N Q(w; x_n) = \mathbb{E}_{\mathbf{x}_i^{\text{emp}}} Q(w; \mathbf{x}_i^{\text{emp}}) = \mathbb{E}_{\mathbf{x}^{\text{emp}}} Q(w; \mathbf{x}^{\text{emp}}) \quad (83)$$

In other words, the empirical risk minimization problem (78) is a special case of the expected risk minimization problem (74), where we define a new random variable  $\mathbf{x}^{\text{emp}} \triangleq \mathbf{x}_n$ , which follows a uniform empirical distribution over the sample batch of data  $\{x_n\}_{n=1}^N$ . Note that for any finite  $N$ , there will be a difference between  $\mathbf{x}$ , the underlying random variable that led to the batch of samples  $\{x_n\}_{n=1}^N$ , and  $\mathbf{x}^{\text{emp}}$ , which is sampled from  $\{x_n\}_{n=1}^N$ .



# Stochastic Gradient Approximations

Given a risk function  $J(w)$  we denote by  $\widehat{\nabla J}(w)$  a stochastic gradient approximation for the true gradient  $\nabla J(w)$  if  $\widehat{\nabla J}(w)$  can be evaluated solely from data available to the algorithm and

$$\widehat{\nabla J}(w) \approx \nabla J(w) \quad (84)$$

in some sense. On occasion, we may wish to make precise the exact data used to compute  $\widehat{\nabla J}(w)$ , in which case we may for example write:

$$\widehat{\nabla J}(w; \mathbf{x}_i) \approx \nabla J(w) \quad (85)$$

to emphasize that  $\widehat{\nabla J}(w; \mathbf{x}_i)$  is constructed based on a realization of the random variable  $\mathbf{x}_i$ . As a general rule, we will employ  $\widehat{\nabla J}(w)$  whenever the data used to construct the gradient approximation is either clear from the context (to lighten the notation), or irrelevant to the discussion (to keep the results general).

## Gradient Noise

Motivated by the discussion so far, we will now study generic expected risk minimization problems of the form:

$$\arg \min_{w \in \mathbb{R}^M} J(w) \triangleq \arg \min_{w \in \mathbb{R}^M} \mathbb{E}_{\mathbf{x}} Q(w; \mathbf{x}) \quad (86)$$

with the understanding that this formulation covers empirical risk minimization problems as laid out in the previous example. We will consider a stochastic gradient algorithm with a generic gradient approximation:

$$\mathbf{w}_i = \mathbf{w}_{i-1} - \mu \widehat{\nabla J}(\mathbf{w}_{i-1}) \quad (87)$$

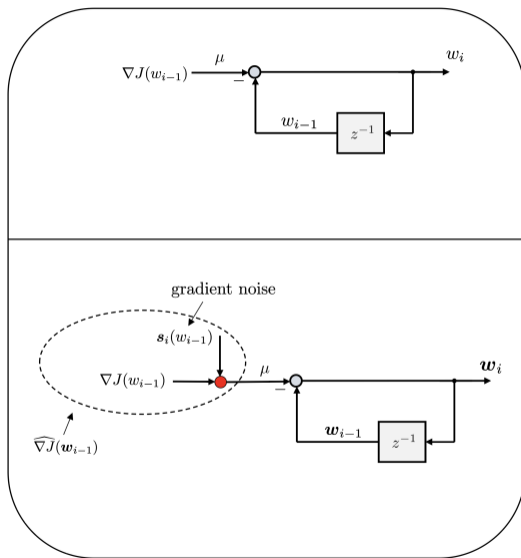
It is instructive to reformulate the recursion as (see Fig. 59):

$$\mathbf{w}_i = \mathbf{w}_{i-1} - \mu \nabla J(\mathbf{w}_{i-1}) - \mu \mathbf{s}_i(\mathbf{w}_{i-1}) \quad (88)$$

in terms of the *gradient noise* term defined by

$$\mathbf{s}_i(\mathbf{w}_{i-1}) = \widehat{\nabla J}(\mathbf{w}_{i-1}) - \nabla J(\mathbf{w}_{i-1}) \quad (89)$$

# Stochastic Gradient Algorithms as Perturbed Systems



## Gradient Noise Modeling Conditions

The gradient approximations  $\widehat{\nabla J}(\mathbf{w}_{i-1})$  are unbiased conditioned on the iterate  $\mathbf{w}_{i-1}$ . Specifically:

$$\mathbb{E} \left\{ \widehat{\nabla J}(\mathbf{w}_{i-1}) \mid \mathbf{w}_{i-1} \right\} = \nabla J(\mathbf{w}_{i-1}) \iff \mathbb{E} \{ \mathbf{s}_i(\mathbf{w}_{i-1}) \mid \mathbf{w}_{i-1} \} = 0 \quad (90)$$

Unbiased gradient estimates ensure that the gradient approximation  $\widehat{\nabla J}(\mathbf{w}_{i-1})$  is, on average, a good estimate for the descent direction. Unbiased estimates are good estimates, but only if their variance is bounded. We hence impose an additional condition:

$$\mathbb{E} \left\{ \|\mathbf{s}_i(\mathbf{w}_{i-1})\|^2 \mid \mathbf{w}_{i-1} \right\} \leq \alpha^2 \|\nabla J(\mathbf{w}_{i-1})\|^2 + \beta^2 \|w^o - \mathbf{w}_{i-1}\|^2 + \gamma^2 (J(\mathbf{w}_{i-1}) - J^o) + \sigma^2 \quad (91)$$

This second condition imposes a *relative* bound on the variance of the gradient approximation because we allow the right-hand side of (91) to grow with various measures of suboptimality.

## On the generality of condition (91)

We note that the variance condition (91) is deliberately chosen to be general by allowing for relative components proportional to  $\|\nabla J(\mathbf{w}_{i-1})\|^2$ ,  $\|w^o - \mathbf{w}_{i-1}\|^2$ , and  $J(\mathbf{w}_{i-1}) - J^o$  through the constants  $\alpha$ ,  $\beta$  and  $\gamma$ , respectively. It is useful to observe that for objectives with  $\delta$ -Lipschitz gradients, it holds that:

$$\frac{1}{2\delta} \|\nabla J(\mathbf{w}_{i-1})\|^2 \leq J(\mathbf{w}_{i-1}) - J^o \leq \frac{\delta}{2} \|w^o - \mathbf{w}_{i-1}\|^2 \quad (92)$$

As such, any gradient approximation that satisfies (91) will also satisfy:

$$\mathbb{E} \left\{ \|\mathbf{s}_i(\mathbf{w}_{i-1})\|^2 \mid \mathbf{w}_{i-1} \right\} \leq \beta^2 \|w^o - \mathbf{w}_{i-1}\|^2 + \sigma^2 \quad (93)$$

for some  $\beta^2$  that is suitably adjusted. For this reason, at several points in our presentation, and in order to lighten the notation, we will set  $\alpha = \gamma = 0$  and work with the simplified condition (93). Nevertheless, allowing for more granular characterization of the relative components in (91) is useful depending on the setting, so we allow for it when it doesn't result in unnecessary complications.

## Elementary Gradient Approximations

Perhaps the most commonly employed gradient approximation is the one obtained from the expression  $\nabla J(w) \triangleq \nabla \{\mathbb{E}_{\mathbf{x}} Q(w; \mathbf{x})\}$  by dropping the expectation and using the loss value at the sample  $\mathbf{x}_i$  available at time  $i$ . This results in:

$$\widehat{\nabla J}^{\text{ele}}(w) \triangleq \nabla Q(w; \mathbf{x}_i) \quad (94)$$

We will be referring to such approximations as *elementary* stochastic gradient approximations. They are the most common and simplest constructions, and we will see that many more elaborate constructions can be “built” from elementary gradient approximations. When we wish to emphasize the fact that the construction is of the elementary type, or to distinguish it from other constructions, we will attach the superscript “ele” to the relevant quantities. For example,  $\widehat{\nabla J}^{\text{ele}}(w)$  refers to the elementary gradient approximation constructed according to (94), while  $s_i^{\text{ele}}(w)$  refers to the resulting gradient noise with parameters  $\alpha_{\text{ele}}^2, \beta_{\text{ele}}^2, \gamma_{\text{ele}}^2, \sigma_{\text{ele}}^2$ . When the context is clear, or we wish to preserve generality of the statements, we will omit “ele”.

## Example: Mean-square error learning

Let us consider the linear regression model:

$$\boldsymbol{\gamma} = \mathbf{h}^T w^o + \mathbf{v} \quad (95)$$

where  $\mathbf{h} \in \mathbb{R}^M$  denotes the regressor and  $\boldsymbol{\gamma}$  denotes the observation. The measurement noise is denoted by  $\mathbf{v}$  and assumed to be zero-mean and independent of  $\mathbf{h}$ . Given observations  $\mathbf{x} \triangleq \text{col}\{\mathbf{h}, \boldsymbol{\gamma}\}$ , we may find the unknown model  $w^o$  by solving:

$$w^o = \arg \min_{w \in \mathbb{R}^M} \frac{1}{2} \mathbb{E} \|\boldsymbol{\gamma} - \mathbf{h}^T w\|^2 \quad (96)$$

Let  $Q(w; \mathbf{x}) = Q(w; \mathbf{h}, \boldsymbol{\gamma}) = \frac{1}{2} \|\boldsymbol{\gamma} - \mathbf{h}^T w\|^2$ . By differentiating, we find:

$$\nabla J(w) = \mathbb{E} \left\{ -\mathbf{h} \left( \boldsymbol{\gamma} - \mathbf{h}^T w \right) \right\} = -r_{\mathbf{h}\boldsymbol{\gamma}} + R_{\mathbf{h}} w \quad (97)$$

where we defined  $r_{\mathbf{h}\boldsymbol{\gamma}} = \mathbb{E} \mathbf{h} \boldsymbol{\gamma}$  and  $R_{\mathbf{h}} \triangleq \mathbb{E} \mathbf{h} \mathbf{h}^T$ .

## Example: Mean-square error learning

We can then construct:

$$\widehat{\nabla J}(w) = \nabla Q(w; \mathbf{x}_i) = -\mathbf{h}_i \left( \gamma_i - \mathbf{h}_i^\top w \right) = -\mathbf{h}_i \gamma_i + \mathbf{h}_i \mathbf{h}_i^\top w \quad (98)$$

so that the gradient noise  $\mathbf{s}_i(w)$  is given by

$$\mathbf{s}_i(w) = -\mathbf{h}_i \gamma_i + \mathbf{h}_i \mathbf{h}_i^\top w + r_{\mathbf{h}\gamma} - R_{\mathbf{h}} w \quad (99)$$

It can be verified that  $\mathbf{s}_i(\mathbf{w}_{i-1})$  is zero-mean and satisfies the variance bound with  $\alpha^2 = \gamma^2 = 0$  and:

$$\beta^2 \triangleq \mathbb{E} \left\| \mathbf{h}_i \mathbf{h}_i^\top - R_{\mathbf{h}} \right\|^2 \quad (100)$$

$$\sigma^2 \triangleq \sigma_v^2 \text{Tr} (R_{\mathbf{h}}) \quad (101)$$



## Example: Logistic Regression

In this case, the labels  $\gamma$  are sampled from a discrete distribution, say  $\gamma = \pm 1$ , while

$$Q(w; \mathbf{h}_i, \gamma_i) \triangleq \ln \left( 1 + e^{-\gamma_i \mathbf{h}_i^\top w} \right) + \frac{\rho}{2} \|w\|^2 \quad (102)$$

$$J(w) \triangleq \mathbb{E} Q(w; \mathbf{h}_i, \gamma_i) = \mathbb{E} \ln \left( 1 + e^{-\gamma_i \mathbf{h}_i^\top w} \right) + \frac{\rho}{2} \|w\|^2 \quad (103)$$

Then, it follows that:

$$\begin{aligned} \nabla Q(w; \mathbf{h}_i, \gamma_i) &= - \frac{\gamma_i \mathbf{h}_i}{1 + e^{\gamma_i \mathbf{h}_i^\top w}} + \rho w \\ \nabla J(w) &= - \mathbb{E} \left\{ \frac{\gamma_i \mathbf{h}_i}{1 + e^{\gamma_i \mathbf{h}_i^\top w}} \right\} + \rho w \end{aligned} \quad (104)$$

Hence, we can construct  $\widehat{\nabla J}(w) = \nabla Q(w; \mathbf{h}_i, \gamma_i)$  and immediately verify the zero-mean condition (90).

## Example: Logistic Regression

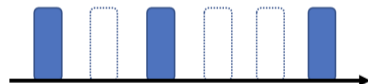
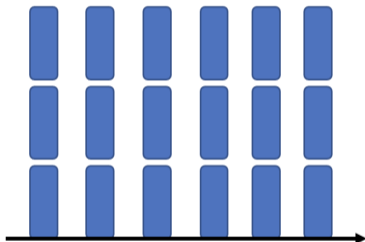
For the variance bound, we can verify that:

$$\mathbb{E} \left\{ \|\mathbf{s}_i(\mathbf{w}_{i-1})\|^2 \mid \mathbf{w}_{i-1} \right\} \leq \mathbb{E} \|\mathbf{h}_i\|^2 = \text{Tr}(R_{\mathbf{h}}) \quad (105)$$

where we defined  $R_{\mathbf{h}} = \mathbb{E} \mathbf{h} \mathbf{h}^{\top}$ . Hence, we know that the elementary logistic regression recursion satisfies the gradient noise conditions with  $\alpha^2 = \beta^2 = \gamma^2 = 0$  and  $\sigma^2 = \text{Tr}(R_{\mathbf{h}})$ .

## Variants of Stochastic Gradient Descent

- We might receive multiple samples at a time, or sometimes receive no samples at all.
- The stochastic approximation framework is flexible to cover these variants.



$$\mathbf{w}_{k,i} = \mathbf{w}_{k,i-1} - \frac{\mu}{3} \sum_{b=1}^3 \nabla Q(\mathbf{w}_{k,i-1}; \mathbf{x}_{k,i,b})$$

$$\mathbf{w}_{k,i} = \begin{cases} \mathbf{w}_{k,i-1} - \frac{\mu}{\pi} \nabla Q(\mathbf{w}_{k,i-1}; \mathbf{x}_{k,i}) & \text{w.p. } \pi, \\ \mathbf{w}_{k,i-1} & \text{w.p. } 1 - \pi. \end{cases}$$

## Mini-Batch SGD

Suppose that instead of being provided with a single data point  $\mathbf{x}_i$  at any given time  $i$ , we have access to  $B$  independent samples  $\{\mathbf{x}_{b,i}\}_{b=1}^B$ . We can then compute:

$$\widehat{\nabla J}(\mathbf{w}_{i-1}; \{\mathbf{x}_{b,i}\}_{b=1}^B) \triangleq \frac{1}{B} \sum_{b=1}^B \nabla Q(\mathbf{w}_{i-1}; \mathbf{x}_{b,i}) \quad (106)$$

It is straightforward to verify that:

$$\mathbb{E} \left\{ \widehat{\nabla J}(\mathbf{w}_{i-1}; \{\mathbf{x}_{b,i}\}_{b=1}^B) \mid \mathbf{w}_{i-1} \right\} = \nabla J(\mathbf{w}_{i-1}) \quad (107)$$

Using the independence of samples, we have for the resulting gradient noise process:

$$\mathbb{E} \left\{ \|\mathbf{s}_i(\mathbf{w}_{i-1})\|^2 \mid \mathbf{w}_{i-1} \right\} \leq \alpha_B^2 \|\nabla J(\mathbf{w}_{i-1})\|^2 + \beta_B^2 \|\mathbf{w}^o - \mathbf{w}_{i-1}\|^2 + \gamma_B^2 (J(\mathbf{w}_{i-1}) - J^o) + \sigma_B^2 \quad (108)$$

with

$$\alpha_B^2 \triangleq \frac{\alpha_{\text{ele}}^2}{B}, \quad \beta_B^2 \triangleq \frac{\beta_{\text{ele}}^2}{B}, \quad \gamma_B^2 \triangleq \frac{\gamma_{\text{ele}}^2}{B}, \quad \sigma_B^2 \triangleq \frac{\sigma_{\text{ele}}^2}{B} \quad (109)$$

## Stochastic Gradient with Missing Samples

We define an i.i.d. Bernoulli random indicator variable:

$$\mathbb{I}_i = \begin{cases} 1, & \text{with probability } \pi, \\ 0, & \text{otherwise.} \end{cases} \quad (110)$$

where  $\mathbb{I}_i$  indicates whether we receive a sample  $\mathbf{x}_i$  at time  $i$  or not. We can then construct:

$$\widehat{\nabla J}(w) \triangleq \begin{cases} \frac{1}{\pi} \nabla Q(w; \mathbf{x}_i), & \text{if } \mathbb{I}_i = 1, \\ 0, & \text{otherwise.} \end{cases} \quad (111)$$

Note that in the case when data is available, the gradient approximation  $\frac{1}{\pi} \nabla Q(w; \mathbf{x}_i)$  is scaled by  $\frac{1}{\pi}$ . And since  $0 < \pi \leq 1$ , this results in larger steps being taken when data is available to compensate for the fact that with probability  $1 - \pi$ , no update occurs. This scaling ensures that we continue to have an unbiased approximation:

$$\mathbb{E} \widehat{\nabla J}(w) = \mathbb{E} \{ \nabla Q(w; \mathbf{x}_i) \} = \nabla J(w) \quad (112)$$

# Stochastic Gradient with Missing Samples

For the gradient noise variance, we have:

$$\begin{aligned} & \mathbb{E} \left\{ \|\mathbf{s}_i(\mathbf{w}_{i-1})\|^2 \mid \mathbf{w}_{i-1} \right\} \\ & \leq \alpha_{\text{asy}}^2 \|\nabla J(\mathbf{w}_{i-1})\|^2 + \beta_{\text{asy}}^2 \|\mathbf{w}^o - \mathbf{w}_{i-1}\|^2 + \gamma_{\text{asy}}^2 (J(\mathbf{w}_{i-1}) - J^o) + \sigma_{\text{asy}}^2 \end{aligned} \quad (113)$$

with:

$$\alpha_{\text{asy}}^2 \triangleq \frac{\alpha_{\text{ele}}^2 + 1 - \pi}{\pi}; \quad \beta_{\text{asy}}^2 \triangleq \frac{\beta_{\text{ele}}^2}{\pi}; \quad \gamma_{\text{asy}}^2 \triangleq \frac{\gamma_{\text{ele}}^2}{\pi}; \quad \sigma_{\text{asy}}^2 \triangleq \frac{\sigma_{\text{ele}}^2}{\pi} \quad (114)$$

# Convergence Guarantee of Stochastic Gradient Algorithms

## Mean-squared deviation of stochastic gradient algorithms

Let the risk function  $J(w)$  be  $\nu$ -strongly convex with  $\delta$ -Lipschitz gradients. Suppose that we employ a gradient approximation  $\widehat{\nabla}J(\cdot)$  satisfying conditions (90) and (91) with constants  $\alpha^2, \beta^2, \gamma^2, \sigma^2 \geq 0$ . Then, the mean-squared deviation  $\mathbb{E}\tilde{\mathbf{w}}_i \triangleq \mathbb{E}\|w^o - \mathbf{w}_i\|^2$  of the iterates generated by the stochastic gradient algorithm (87) satisfies:

$$\mathbb{E}\|\tilde{\mathbf{w}}_i\|^2 \leq \lambda \mathbb{E}\|\tilde{\mathbf{w}}_{i-1}\|^2 + \mu^2 \sigma^2 \quad (115)$$

where we defined

$$\lambda \triangleq 1 - 2\mu\nu + \mu^2 \left( (1 + \alpha^2)\delta^2 + \beta^2 + \gamma^2 \frac{\delta}{2} \right) \quad (116)$$

# Convergence Guarantee of Stochastic Gradient Algorithms

## Mean-squared deviation of stochastic gradient algorithms (2)

Moreover, for sufficiently small step-sizes:

$$\mu \leq \frac{\nu}{(1 + \alpha^2)\delta^2 + \beta^2 + \gamma^2 \frac{\delta}{2}} \quad (117)$$

it holds that  $\lambda < 1$  and we can iterate (115) to find:

$$\mathbb{E}\|\tilde{\mathbf{w}}_i\|^2 \leq \lambda^i \|\tilde{\mathbf{w}}_0\|^2 + \frac{\mu\sigma^2}{\nu} \quad (118)$$

For the excess risk, we obtain:

$$\mathbb{E}J(\mathbf{w}_i) - J(w^o) \leq \lambda^i \left( \frac{\delta \|\tilde{\mathbf{w}}_0\|^2}{2} \right) + \frac{\mu\delta\sigma^2}{2\nu} \quad (119)$$



## Transient and Steady-State Error

We observe that the theorem quantifies both the transient and steady-state behavior of the mean-squared deviation  $\mathbb{E}\|\tilde{\mathbf{w}}_i\|^2$  and the excess-risk  $\mathbb{E}J(\mathbf{w}_i) - J(w^o)$ . The steady-state component of an error refers to the part of the error which remains as  $i \rightarrow \infty$ . For the MSD in (118), this corresponds to:

$$\limsup_{i \rightarrow \infty} \mathbb{E}\|\tilde{\mathbf{w}}_i\|^2 = \frac{\mu\sigma^2}{\nu} = O(\mu) \quad (120)$$

while the transient component corresponds to  $\lambda^i \|\tilde{\mathbf{w}}_0\|^2$ . Similarly, for the excess risk, we find in steady-state:

$$\limsup_{i \rightarrow \infty} \mathbb{E}J(\mathbf{w}_i) - J(w^o) = \frac{\mu\delta\sigma^2}{2\nu} = O(\mu) \quad (121)$$

and identify the transient component as  $\lambda^i \left( \frac{\delta\|\tilde{\mathbf{w}}_0\|^2}{2} \right)$ .

## Alternative Metrics: Iteration Complexity

The iteration complexity of an algorithm is defined as the number of iterations needed for its error to fall below a specified threshold  $\epsilon$ . We can readily obtain iteration complexities from the characterizations (118) and (119) of the MSD and excess risk respectively. We will illustrate this for the MSD — the argument for the excess risk is analogous. Our objective is to determine the number of iterations  $i^o$  needed so that:

$$\mathbb{E}\|\tilde{\mathbf{w}}_{i^o}\|^2 \leq \epsilon \quad (122)$$

We can find the IR from (118) by appropriately choosing the step-size  $\mu$ :

$$i^o \geq \frac{2\sigma^2}{\epsilon} \ln \left( \frac{1}{2\epsilon\|\tilde{\mathbf{w}}_0\|^2} \right) \quad (123)$$

## Alternative Metrics: Expected Regret

A second commonly encountered metric for the performance of online learning algorithms is the expected regret, defined the cumulative excess risk encountered up to time  $T$ :

$$\text{Regret}_T \triangleq \sum_{i=1}^T \mathbb{E}J(\mathbf{w}_i) - J(w^o) \quad (124)$$

By suitably choosing the step-size, we can find the expected regret from (119) to be:

$$\sum_{i=1}^T (\mathbb{E}J(\mathbf{w}_i) - J(w^o)) = O\left(\sqrt{T}\right) \quad (125)$$

## Case Study: Mean-Square Error

We illustrate the flexibility of the theorem in combination with the stochastic gradient approximation framework in this chapter by developing a number of variants of the stochastic gradient algorithm along with the corresponding performance bounds. First, let us recall the mean-square error problem, which is given by:

$$J(w) = \frac{1}{2} \mathbb{E} \|\gamma - \mathbf{h}^\top w\|^2 \quad (126)$$

Given a single pair of observations  $\{\gamma_i, \mathbf{h}_i\}$ , we can construct the elementary gradient approximation:

$$\widehat{\nabla J}^{\text{ele}}(\mathbf{w}_{i-1}) = \widehat{\nabla Q}(\mathbf{w}_{i-1}; \mathbf{h}_i, \gamma_i) = -\mathbf{h}_i \left( \gamma_i - \mathbf{h}_i^\top \mathbf{w}_{i-1} \right) \quad (127)$$

The resulting stochastic gradient algorithm amounts to:

$$\mathbf{w}_i^{\text{ele}} = \mathbf{w}_{i-1}^{\text{ele}} - \mu \widehat{\nabla J}^{\text{ele}}(\mathbf{w}_{i-1}^{\text{ele}}) = \mathbf{w}_{i-1}^{\text{ele}} + \mu \mathbf{h}_i \left( \gamma_i - \mathbf{h}_i^\top \mathbf{w}_{i-1}^{\text{ele}} \right) \quad (128)$$

## Elementary Gradient Approximations

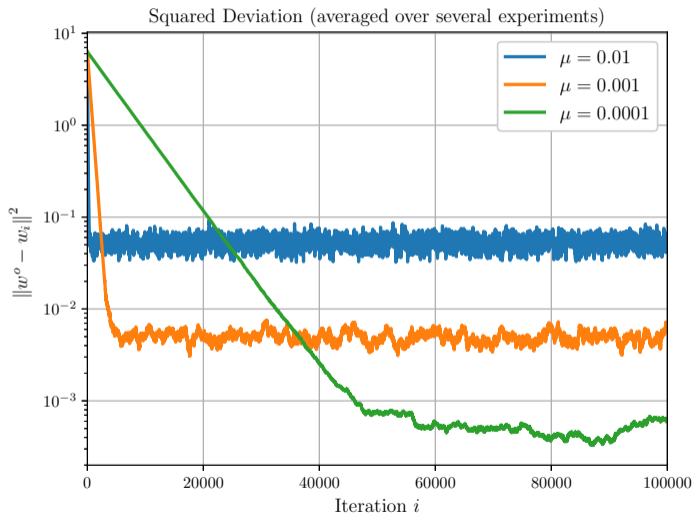
We already computed the resulting gradient noise constants in a previous example to be  $\beta^2 = \mathbb{E} \left\| \mathbf{h}_i \mathbf{h}_i^\top - R_{\mathbf{h}} \right\|^2$  and  $\sigma^2 = \sigma_v^2 \text{Tr} (R_{\mathbf{h}})$  with  $\alpha^2 = \gamma^2 = 0$ . To find the strong-convexity and Lipschitz parameters  $\nu$  and  $\delta$ , we can differentiate (126) twice and find  $\nu = \lambda_{\min} (R_{\mathbf{h}})$  and  $\delta = \lambda_{\max} (R_{\mathbf{h}})$ . From our convergence bounds, we can then conclude directly that:

$$\mathbb{E} \left\| w^o - \mathbf{w}_i^{\text{ele}} \right\|^2 \leq \lambda_{\text{ele}}^i \|\tilde{w}_0\|^2 + \frac{\mu \sigma_v^2 \text{Tr} (R_{\mathbf{h}})}{\lambda_{\min} (R_{\mathbf{h}})} \quad (129)$$

where:

$$\lambda_{\text{ele}} \triangleq 1 - 2\mu \lambda_{\min} (R_{\mathbf{h}}) + \mu^2 \left( \mathbb{E} \left\| \mathbf{h}_i \mathbf{h}_i^\top - R_{\mathbf{h}} \right\|^2 + \lambda_{\max} (R_{\mathbf{h}}) \right) \quad (130)$$

# Learning Dynamics of Elementary Stochastic Gradient



## Mini-Batch Gradient Approximations

If we are instead provide with a mini-batch of observations  $\{\mathbf{h}_{b,i}, \gamma_{b,i}\}_{b=1}^B$  at time  $i$ , we can construct:

$$\widehat{\nabla J}^B(\mathbf{w}_{i-1}) = \frac{1}{B} \sum_{b=1}^B \nabla Q(\mathbf{w}_{i-1}; \mathbf{h}_{b,i}, \gamma_{b,i}) = -\mathbf{h}_{b,i} \left( \gamma_{b,i} - \mathbf{h}_{b,i}^\top \mathbf{w}_{i-1} \right) \quad (131)$$

The resulting mini-batch stochastic gradient algorithm amounts to:

$$\mathbf{w}_i^B = \mathbf{w}_{i-1}^B - \mu \widehat{\nabla J}^B(\mathbf{w}_{i-1}^B) = \mathbf{w}_{i-1}^B + \frac{\mu}{B} \sum_{b=1}^B \mathbf{h}_{b,i} \left( \gamma_{b,i} - \mathbf{h}_{b,i}^\top \mathbf{w}_{i-1}^B \right) \quad (132)$$

We know from (109) that:

$$\beta_B^2 = \frac{\beta^2}{B} = \frac{\mathbb{E} \|\mathbf{h}_i \mathbf{h}_i^\top - R_h\|^2}{B} \quad (133)$$

$$\sigma_B^2 = \frac{\sigma^2}{B} = \frac{\sigma_v^2 \text{Tr}(R_h)}{B} \quad (134)$$

# Mini-Batch Gradient Approximations

We can then conclude again from (118):

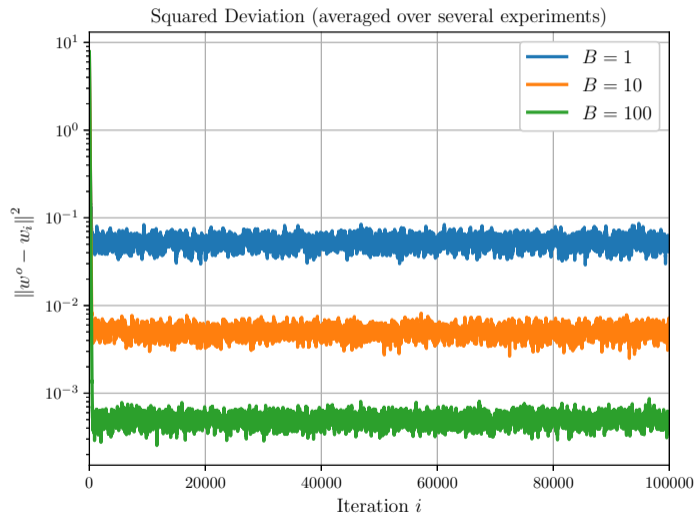
$$\mathbb{E} \|w^o - \mathbf{w}_i^B\|^2 \leq \lambda_B^i \|\tilde{w}_0\|^2 + \frac{\mu \sigma_v^2 \text{Tr}(R_{\mathbf{h}})}{B \lambda_{\min}(R_{\mathbf{h}})} \quad (135)$$

where:

$$\lambda_B \triangleq 1 - 2\mu \lambda_{\min}(R_{\mathbf{h}}) + \mu^2 \left( \frac{\mathbb{E} \|\mathbf{h}_i \mathbf{h}_i^{\top} - R_{\mathbf{h}}\|^2}{B} + \lambda_{\max}(R_{\mathbf{h}}) \right) \quad (136)$$



# Learning Dynamics of Mini-Batch Gradient Approximations



## Stochastic Gradients with Missing Samples

If we are provided with a pair of observations  $\{\mathbf{h}_i, \gamma_i\}$  only with probability  $0 < \pi \leq 1$ , we can construct:

$$\widehat{\nabla J}^{asy}(\mathbf{w}_{i-1}) \triangleq \begin{cases} \frac{-1}{\pi} \mathbf{h}_i (\gamma_i - \mathbf{h}_i^\top \mathbf{w}_{i-1}), & \text{with prob. } \pi, \\ 0, & \text{otherwise.} \end{cases} \quad (137)$$

Then, the resulting stochastic gradient algorithm is given by:

$$\mathbf{w}_i^{asy} = \begin{cases} \mathbf{w}_{i-1}^{asy} + \frac{1}{\pi} \mathbf{h}_i (\gamma_i - \mathbf{h}_i^\top \mathbf{w}_{i-1}^{asy}), & \text{with prob. } \pi, \\ \mathbf{w}_{i-1}^{asy}, & \text{otherwise.} \end{cases} \quad (138)$$

with gradient noise constants modified from those of the elementary approximation via (114):

$$\beta_{asy}^2 = \frac{\beta^2}{\pi} = \frac{\mathbb{E} \|\mathbf{h}_i \mathbf{h}_i^\top - R_h\|^2}{\pi} \quad (139)$$

$$\sigma_{asy}^2 = \frac{\sigma^2}{\pi} = \frac{\sigma_v^2 \text{Tr}(R_h)}{\pi} \quad (140)$$

# Stochastic Gradients with Missing Samples

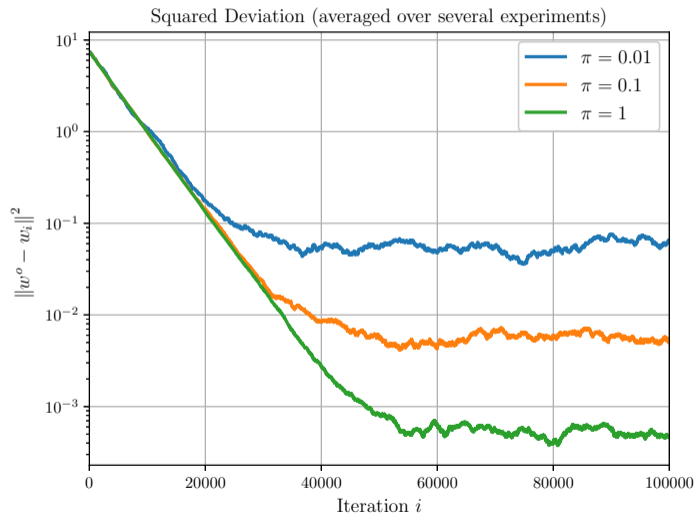
Hence, we find from (118):

$$\mathbb{E}\|w^o - \mathbf{w}_i^{\text{asy}}\|^2 \leq \lambda_{\text{asy}}^i \|\tilde{w}_0\|^2 + \frac{\mu \sigma_v^2 \text{Tr}(R_h)}{\pi \lambda_{\min}(R_h)} \quad (141)$$

where:

$$\lambda_{\text{asy}} \triangleq 1 - 2\mu \lambda_{\min}(R_h) + \mu^2 \left( \frac{\mathbb{E}\|\mathbf{h}_i \mathbf{h}_i^\top - R_h\|^2}{\pi} + \lambda_{\max}(R_h) \right) \quad (142)$$

# Performance of Stochastic Gradients with Missing Samples



# Conclusion

- We saw how:
  - ▶ Multi-agent learning problems naturally lead to individual and aggregate risk minimization problems.
  - ▶ How stochastic gradient algorithms can solve these risk minimization problems with performance generally limited by a trade-off between quality of gradient approximation, rate of convergence and steady-state error.
- Moving forward, we will develop multi-agent learning algorithms that allow agents to collaboratively learn better/faster/more efficiently than any individual agent.
  - ▶ Next lecture we will focus on fusion-center based approaches.
  - ▶ From tomorrow, we will look at decentralized architectures.

# References and Further Reading

- General references and surveys:
  - ▶ A. H. Sayed, *Inference and Learning from Data*, Cambridge University Press, 2022.